

# EECE 5639 Computer Vision I

---

## Lecture 22

**SFM; Object Recognition**

**Project 4 is out. Now Due April 16.**

**Hw 5 is out. Now Due April 19**

## Next Class

**Object Recognition**

# SFM by Factorization

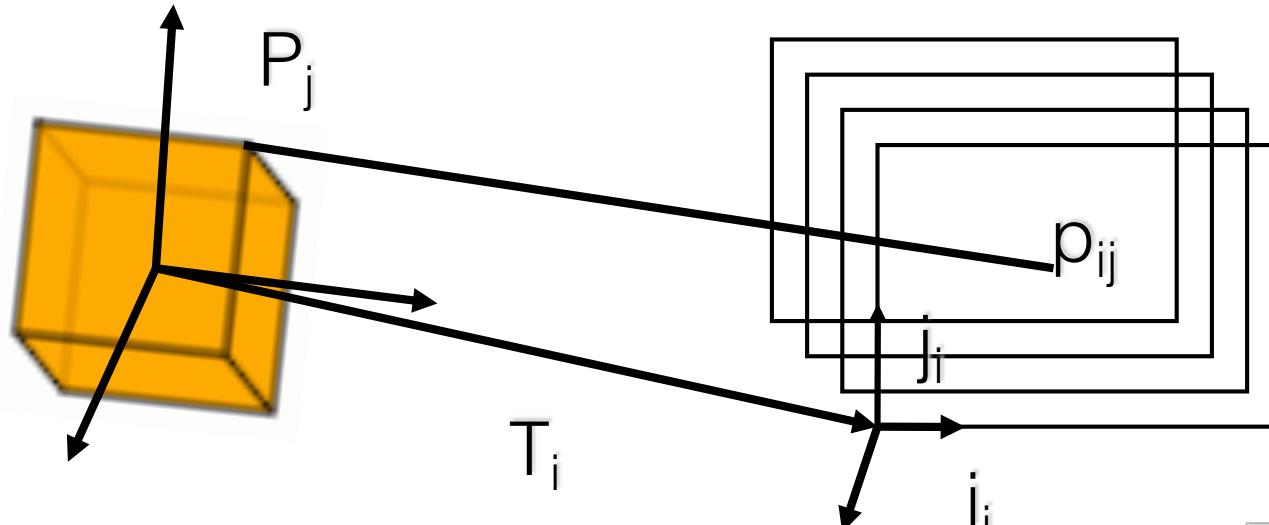
---

Consider a **FIXED** camera and a **MOVING** object

Track n features across m frames.

We want to find the structure and the motion of the object.

# Orthographic Projection

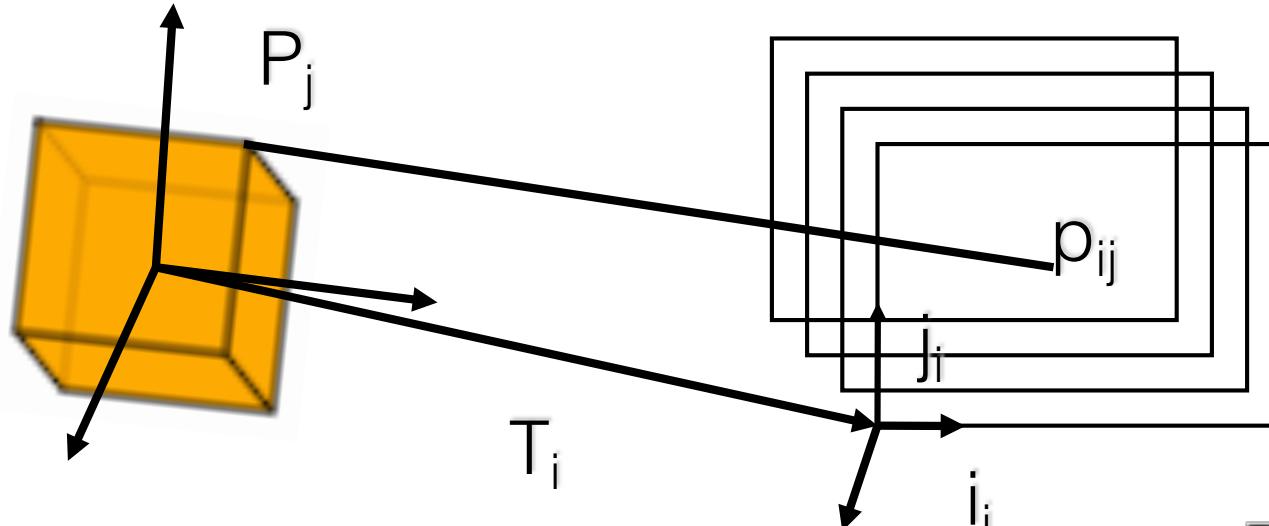


$$p_{ij} = \begin{matrix} R_j \\ 2 \times 3 \end{matrix} \begin{matrix} P_j \\ 3 \times 1 \end{matrix} + \begin{matrix} T_i \\ 2 \times 1 \end{matrix}$$

$$\begin{bmatrix} p_{11} & p_{12} & \dots & p_{1n} \\ p_{21} & p_{22} & \dots & p_{2n} \\ \vdots & \vdots & \vdots & \vdots \\ p_{m1} & p_{m2} & \dots & p_{mn} \end{bmatrix} = \begin{bmatrix} R_1 & T_1 \\ R_2 & T_1 \\ \vdots & \vdots \\ R_m & T_m \end{bmatrix} \begin{bmatrix} P_1 & S & \dots & P_n \\ 1 & 1 & \dots & 1 \end{bmatrix}$$

$$W = MS$$

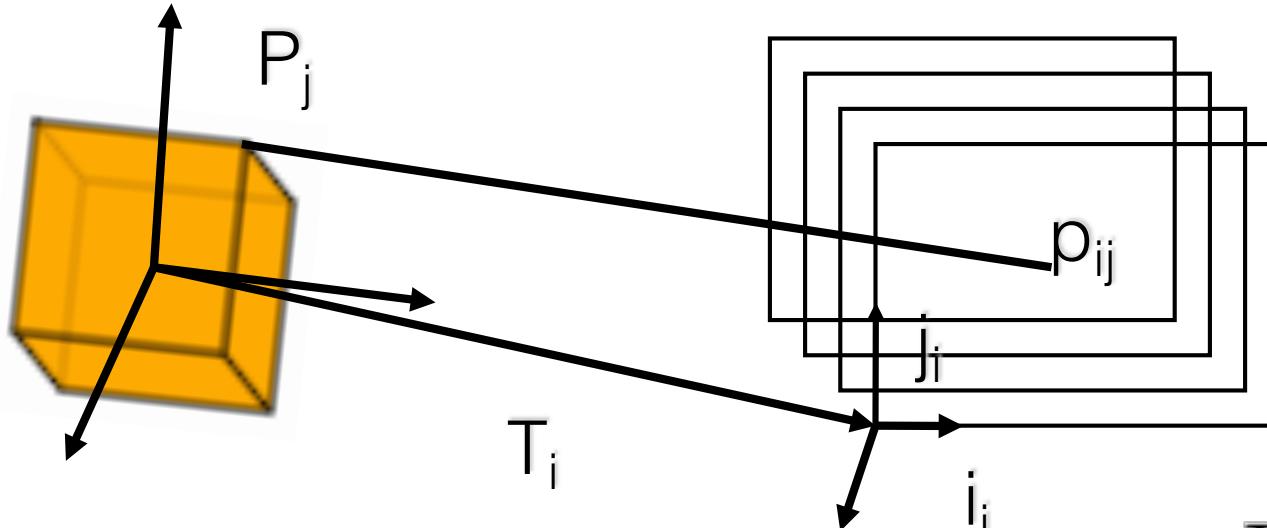
# Orthographic Projection



$$p_{ij} = \begin{matrix} R_i \\ 2 \times 3 \end{matrix} \begin{matrix} P_j \\ 3 \times 1 \end{matrix} + \begin{matrix} T_i \\ 2 \times 1 \end{matrix}$$

One can eliminate  $T$  by setting the origin of both coordinate systems (the object and the image coordinate systems) at the centroids of all the points (3D and each 2D frame).

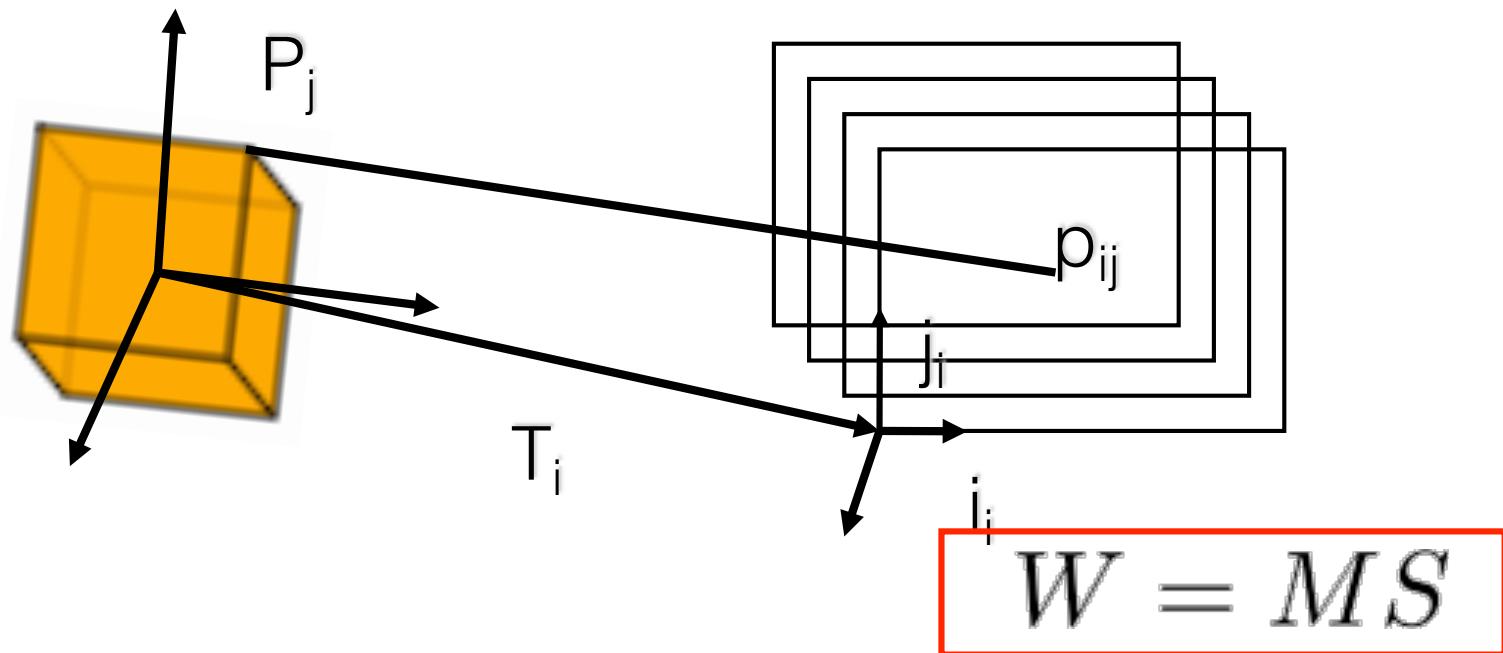
# Orthographic Projection



$$p_{ij} = R_i P_j + T_i$$

$$\begin{bmatrix} \hat{p}_{11} & \hat{p}_{12} & \dots & \hat{p}_{1n} \\ \hat{p}_{21} & \hat{p}_{22} & \dots & \hat{p}_{2n} \\ \vdots & \vdots & \dots & \vdots \\ \hat{p}_{m1} & \hat{p}_{m2} & \dots & \hat{p}_{mn} \end{bmatrix} = \begin{bmatrix} R_1 \\ M \\ R_2 \\ \vdots \\ R_m \end{bmatrix} \begin{bmatrix} \hat{P}_1 & \hat{P}_2 & \dots & \hat{P}_n \end{bmatrix}$$

$W = MS$



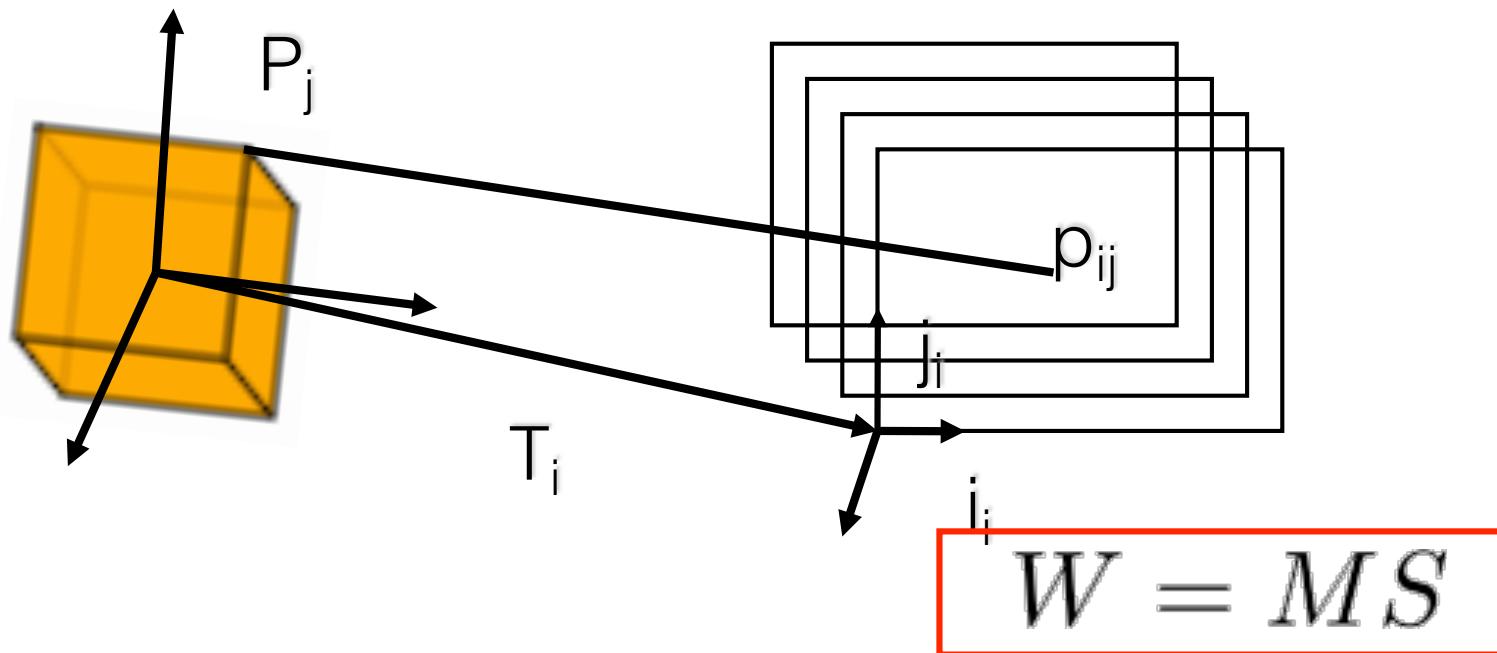
$2m \times n \quad 2m \times 3 \quad 3 \times n$

$W$  is **rank 3** and  $M$  and  $S$  can be found (up to an affine transf.) from the SVD of  $W$

$$W = UDV^T$$

$$M = U D^{\frac{1}{2}} Q \quad S = Q^{-1} D^{\frac{1}{2}} V^T$$

# Finding Q



$$W = UDV^T$$

$$S = Q^{-1} D^{\frac{1}{2}} V^T$$

$$M = U D^{\frac{1}{2}} Q$$

M should be a stack of  
ROTATIONS!

# Finding Q

---

$$M = U D^{\frac{1}{2}} Q$$

M should be a stack of ROTATIONS!

$$m_{i1} \cdot m_{i2} = 0$$

$$|m_{i1}|^2 = 1$$

$$|m_{i2}|^2 = 1$$

$$m_{i1} \times m_{i2} = m_{i3}$$

# Recovering T

---

$$T_i = R_i \begin{bmatrix} \bar{p}_i \\ \alpha \end{bmatrix} \quad \text{with } \alpha \text{ an arbitrary real number}$$

# Algorithm Summary

---

Obtain the registered measurement matrix  $W$  (subtract the mean in each frame)

Compute SVD of  $W = UDV^T$

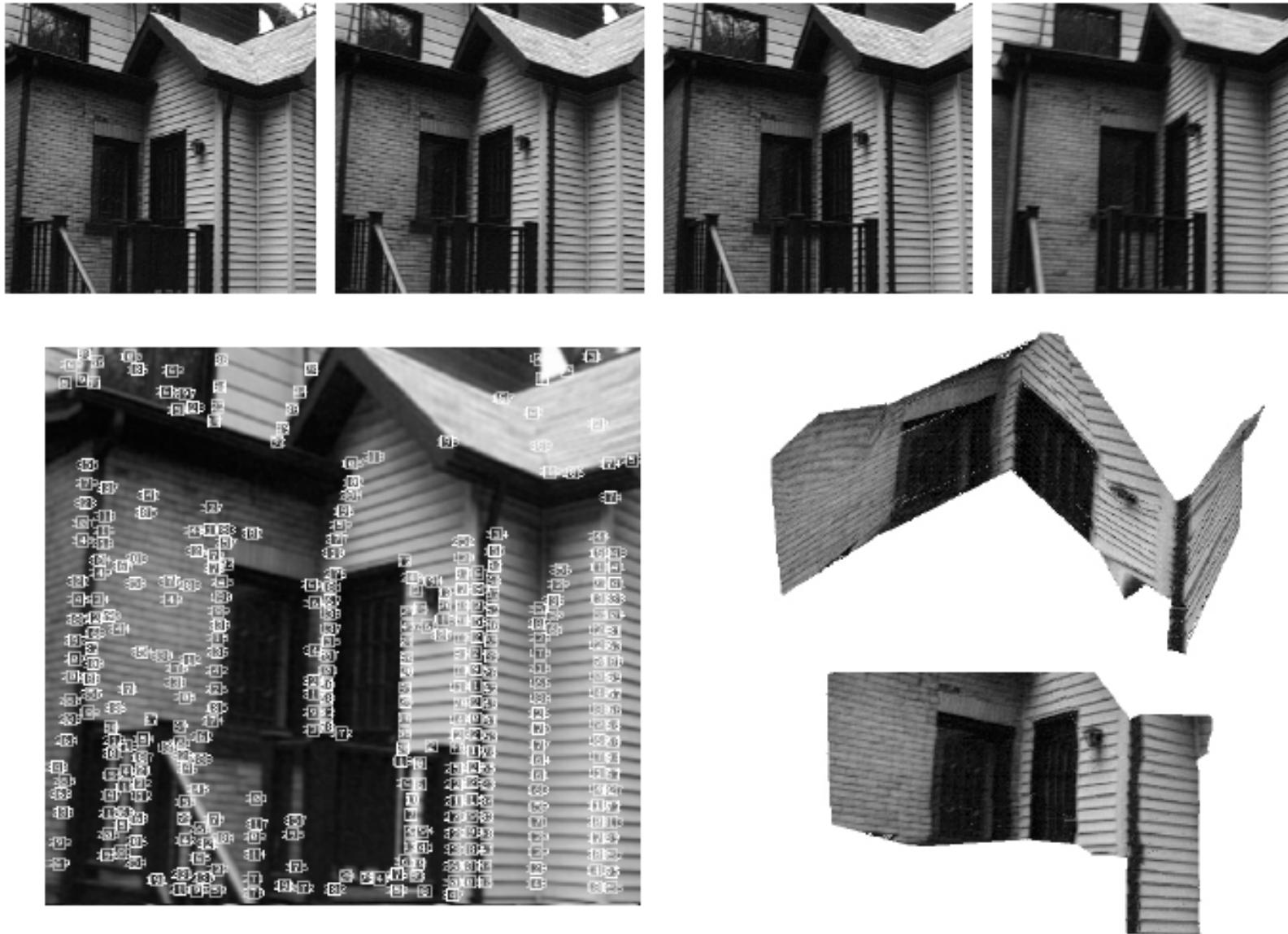
Keep the first 3 eigenvalues/eigenvectors:  $W' = U'D'V'^T$

Define  $M = U'D'^{1/2}$  and  $S = D'^{1/2}V'^T$

Find affinity  $Q$  such that  $MQ$  is a stack of rotations

Compute  $T$

# Reconstruction Results (Tomasi and Kanade, 1992)



Reprinted from Tomasi and Kanade 1992

# Further Factorization work

---

Factorization with uncertainty

Factorization for indep. moving objects

(Irani & Anandan, IJCV'02)

Factorization for dynamic objects

Perspective factorization

(Costeira and Kanade '94)

Factorization with outliers and missing pts.

(Bregler et al. 2000, Brand 2001)

(Sturm & Triggs 1996, Ayazoglu, Sznajer, Camps 2010...)

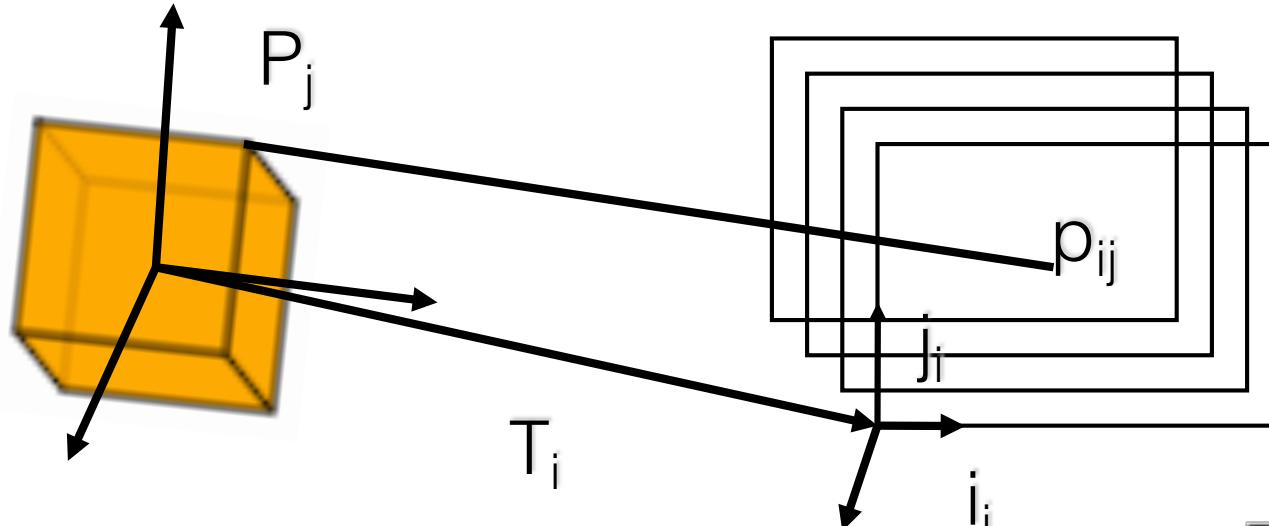
(Jacobs 1997 (affine),  
Martinek and Pajdla 2001,  
Aanaes 2002 (perspective))

# Motion Based Segmentation by Factorization

---

Costeira & Kanade '98

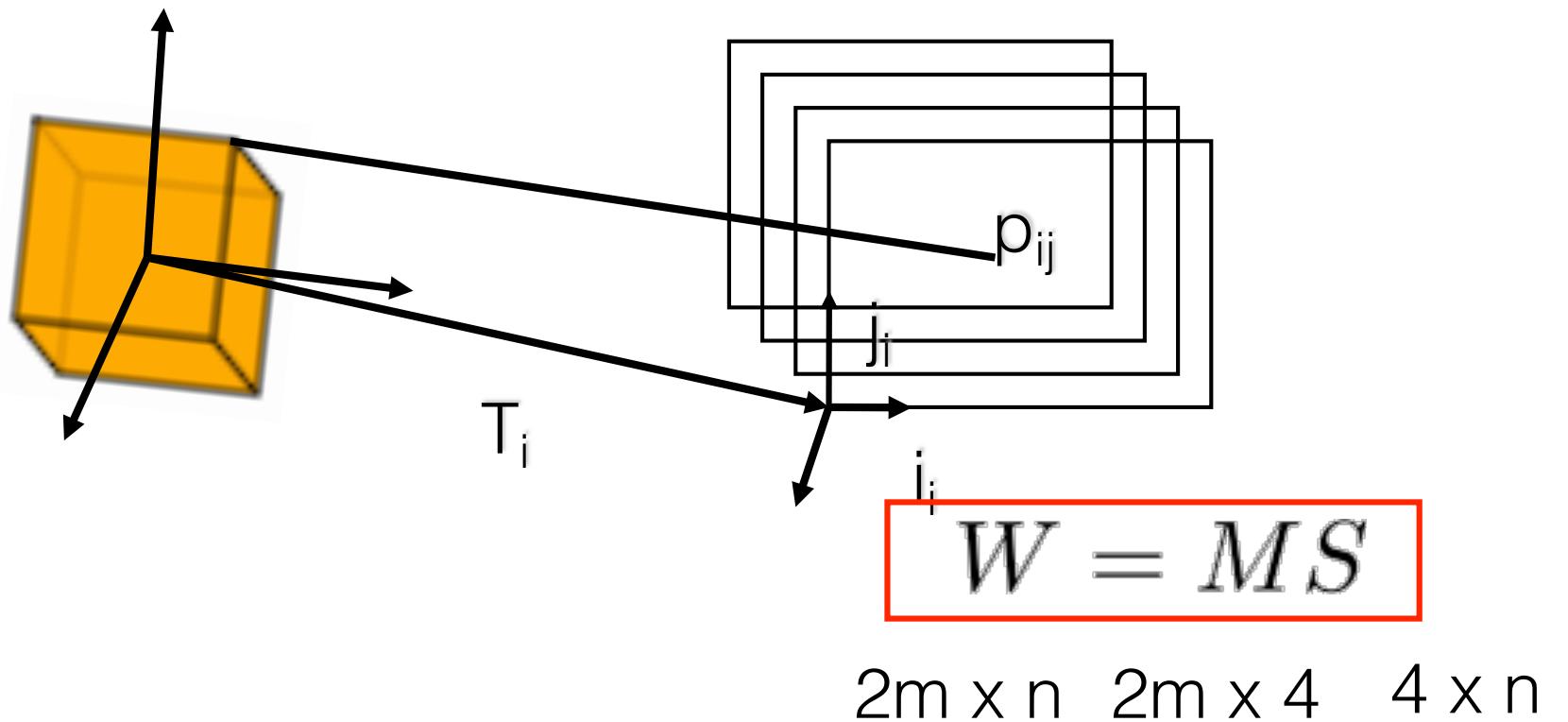
# Motion based Segmentation by Factorization



$$p_{ij} = R_i P_j + T_i$$

$$\begin{bmatrix} p_{11} & p_{12} & \dots & p_{1n} \\ p_{21} & p_{22} & \dots & p_{2n} \\ \vdots & W & \vdots & \vdots \\ p_{m1} & p_{m2} & \dots & p_{mn} \end{bmatrix} = \begin{bmatrix} R_1 & T_1 \\ R_2 & T_1 \\ \vdots & \vdots \\ R_m & T_m \end{bmatrix} \begin{bmatrix} P_1 & P_2 & \dots & P_n \\ 1 & 1 & \dots & 1 \end{bmatrix}$$

$$W = MS$$



W has at most rank 4

# Motion based Segmentation by Factorization

For N (segmented) rigid objects tracked across m frames:

$$\begin{bmatrix} p_{11}^1 & \dots & p_{1n_1}^1 & p_{11}^2 & \dots & p_{1n_2}^2 & \dots & p_{11}^N & \dots & p_{1n_N}^N \\ \vdots & \dots & \vdots & \vdots & \dots & \vdots & \dots & \vdots & \dots & \vdots \\ p_{m1}^1 & \dots & p_{mn_1}^1 & p_{m1}^2 & \dots & p_{mn_2}^2 & \dots & p_{m1}^N & \dots & p_{mn_N}^N \end{bmatrix} = \begin{bmatrix} R_1^1 & T_1^1 & R_1^2 & T_1^2 & \dots & R_1^N & T_1^N \\ R_2^1 & T_2^1 & R_2^2 & T_2^2 & \dots & R_2^N & T_2^N \\ \vdots & \vdots & \vdots & \vdots & \dots & \vdots & \vdots \\ R_m^1 & T_m^1 & R_m^2 & T_m^2 & \dots & R_m^N & T_m^N \end{bmatrix} \begin{bmatrix} P_1^1 & \dots & P_n^1 & 0 & \dots & 0 & \dots & 0 & \dots & 0 \\ 1 & \dots & 1 & 0 & \dots & 0 & \dots & 0 & \dots & 0 \\ 0 & \dots & 0 & P_1^2 & \dots & P_n^2 & \dots & 0 & \dots & 0 \\ 0 & \dots & 0 & 1 & \dots & 1 & \dots & 0 & \dots & 0 \\ 0 & \dots & 0 & 0 & \dots & 0 & \dots & P_1^N & \dots & P_n^N \\ 0 & \dots & 0 & 0 & \dots & 0 & \dots & 1 & \dots & 1 \end{bmatrix}$$

2m x (n<sub>1</sub>+n<sub>2</sub>+...+n<sub>N</sub>)      2m x 4N      4N x (n<sub>1</sub>+n<sub>2</sub>+...+n<sub>N</sub>)

$$W = MS$$

W has at most rank 4N

# Motion based Segmentation by Factorization

For N (segmented) rigid objects tracked across m frames:

$$\begin{bmatrix} p_{11}^1 & \dots & p_{1n_1}^1 \\ \vdots & \dots & \vdots \\ p_{m1}^1 & \dots & p_{mn_1}^1 \\ \vdots & \dots & \vdots \\ p_{11}^N & \dots & p_{1n_N}^N \\ \vdots & \dots & \vdots \\ p_{m1}^N & \dots & p_{mn_N}^N \end{bmatrix} = \begin{bmatrix} R_1^1 & T_1^1 \\ R_2^1 & T_2^1 \\ \vdots & \vdots \\ R_m^1 & T_m^1 \\ R_1^2 & T_1^2 \\ R_2^2 & T_2^2 \\ \vdots & \vdots \\ R_m^2 & T_m^2 \\ \vdots & \vdots \\ R_1^N & T_1^N \\ R_2^N & T_2^N \\ \vdots & \vdots \\ R_m^N & T_m^N \end{bmatrix} \begin{bmatrix} P_1^1 & \dots & P_n^1 & 0 & \dots & 0 & \dots & 0 & 0 \\ 1 & \dots & 1 & 0 & \dots & 0 & \dots & 0 & 0 \\ 0 & \dots & 0 & P_1^2 & \dots & P_n^2 & \dots & 0 & 0 \\ 0 & \dots & 0 & 1 & \dots & 1 & \dots & 0 & 0 \\ 0 & \dots & 0 & 0 & \dots & 0 & \dots & P_1^N & \dots & P_n^N \\ 0 & \dots & 0 & 0 & \dots & 0 & \dots & 1 & \dots & 1 \end{bmatrix}$$

$$2m \times (n_1 + n_2 + \dots + n_N)$$

$$2m \times 4N$$

$$4N \times (n_1 + n_2 + \dots + n_N)$$

$$W = MS$$

$$W = UDV^T$$

$$Z = V(V^T V)^{-1} V^T = VV^T$$

Is block diagonal

# Can we remove the translation component?

---

No ...

---

Each object has a different centroid, but there is only one centroid per frame.

# Motion based Segmentation by Factorization

For N (unsegmented) rigid objects tracked across m frames:

$$\begin{bmatrix} p_{11} & p_{12} & p_{13} & p_{14} & p_{15} & \dots & p_{1K} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ p_{m1} & p_{m2} & p_{m3} & p_{m4} & p_{m5} & \dots & p_{mK} \end{bmatrix}$$

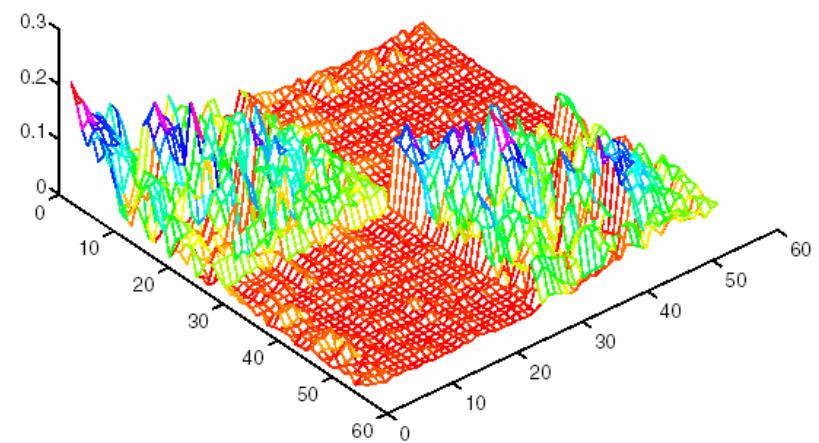
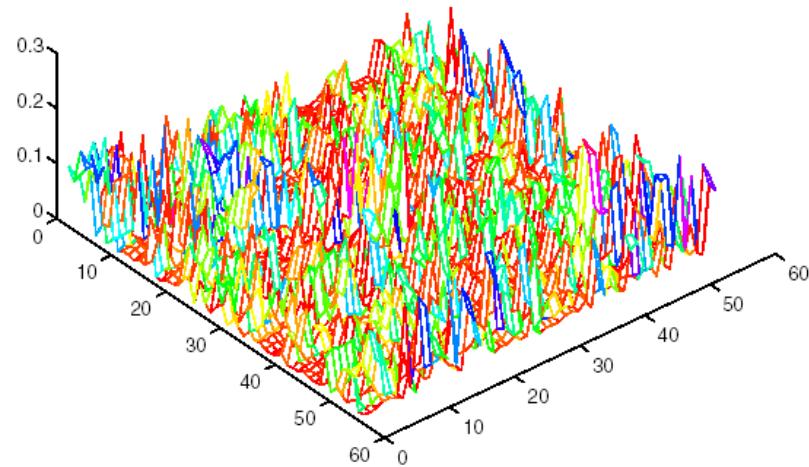
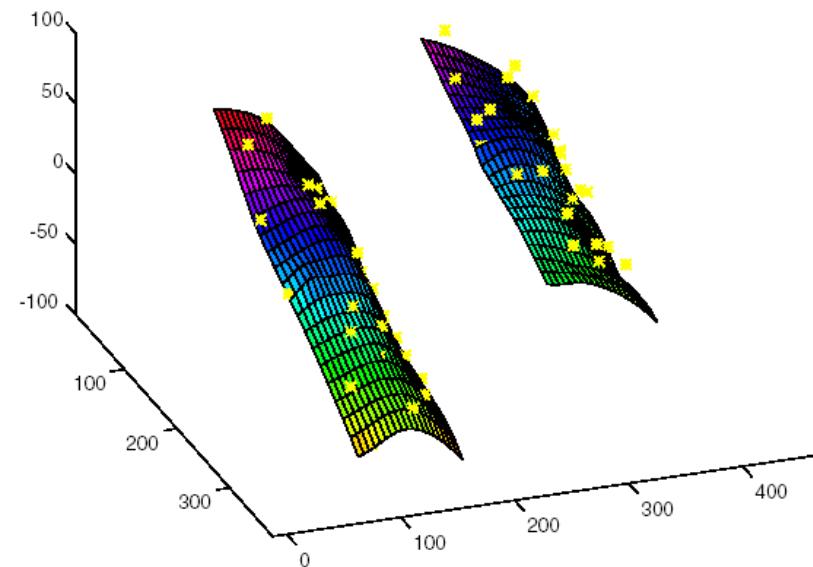
Find SVD of W:  $W = UDV^T$

Find “shape interaction matrix”:

$$Z = VV^T$$

Swap rows and columns of Z until it is block diagonal.

# Multiple indep. moving objects



# Problems:

---

Sensitive to noise, outliers, biased tracking ...

Cannot handle shared motions modes

Ex:

$$M = \begin{bmatrix} M^1 & M^2 & \dots & M^N \end{bmatrix} \text{ With rank 16}$$

How many objects?  $16 = 4N \rightarrow N = 4$

But with shared translation:

$$T^1 = T^2 = \dots = T^N = T$$

M has at most rank  $3N + 1$  !!

$$16 = 3N + 1 \rightarrow N = 5$$

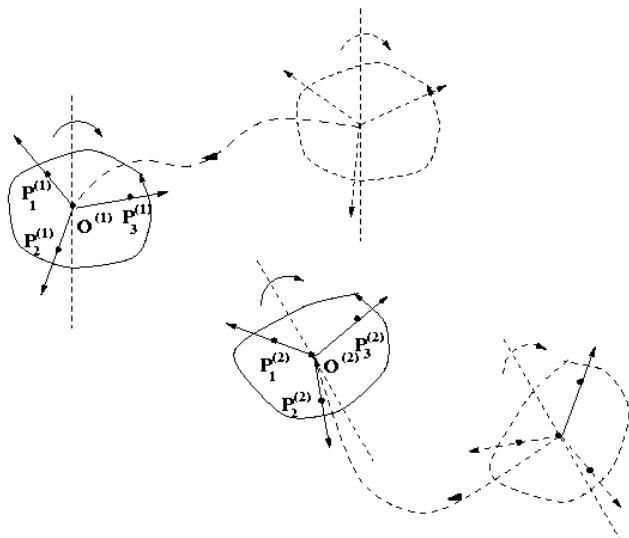
# Motion Segmentation from Dynamics

Lublinerman, Sznajer, Camps '06

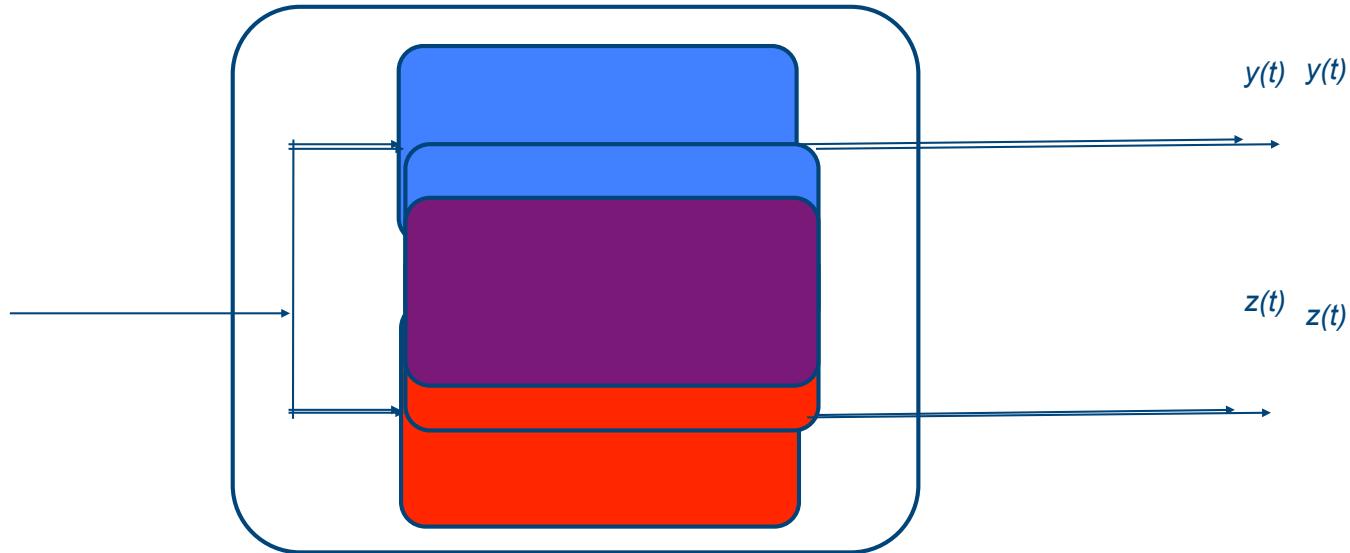
# Observations:

---

- Temporal information is not used: swapping rows in  $W$  does not affect the end result
- The motion of **pairs of points** from a **single object** can be described by the motion of a basis attached to the object but does not carry information about the motion of the other objects.  
**(unobservable states)**
- The motion of **pair of points** from **two different objects** can be described by the motion of **two** basis, each of them attached to one of the objects.



# Dynamic Data Association



A system explaining independent sequences has higher complexity than a system explaining correlated data.

Look for simplest **joint** models.  
*(no need to explicitly find the models!)*

# Main Ideas:

---

Group points according to the complexity of a linear operator required to explain their spatio-temporal evolution.

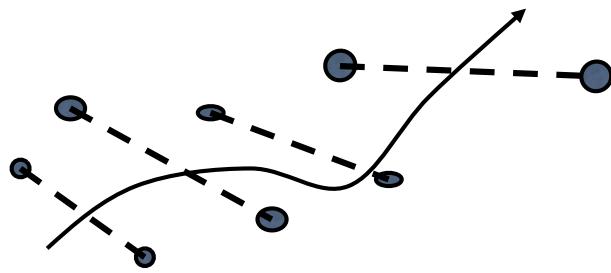
The order of the operator can be estimated by computing the rank of its Hankel matrix.

# Dynamics based Segmentation

$$W = \begin{bmatrix} p_{11} & p_{12} & \dots & p_1 \\ p_{21} & p_{22} & \dots & p_{2n} \\ \vdots & \vdots & \vdots & \vdots \\ p_{m1} & p_{m2} & \dots & p_{mn} \end{bmatrix}$$

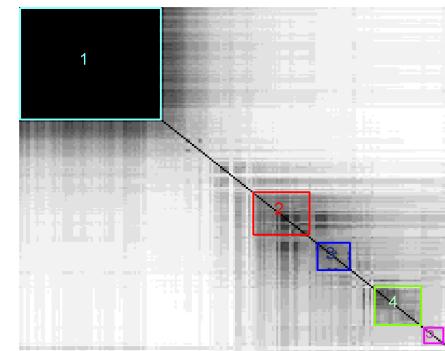
n points, m frames

Relative distance between  $p_i$  and  $p_j$ :

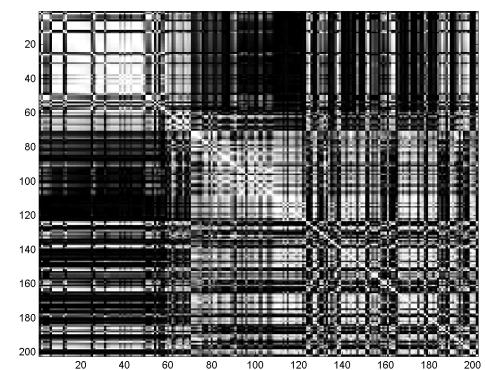
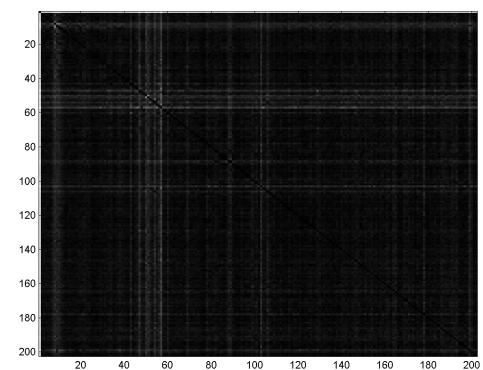
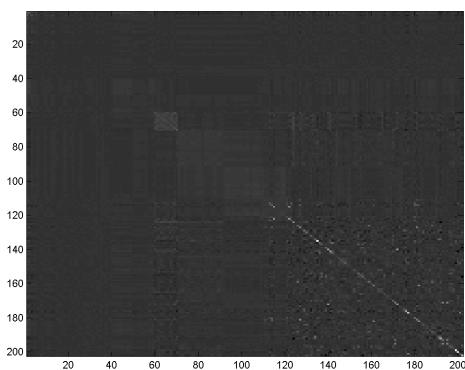


If  $p_i$  and  $p_j$  move together, the order of the dynamics needed to explain their motion has lower order than when they move independently.

$$d_t^{i,j} = p_{ti} - p_{tj}$$
$$H^{i,j} = \begin{bmatrix} d_1^{i,j} & d_2^{i,j} & \dots & d_{m/2}^{i,j} \\ d_2^{i,j} & d_3^{i,j} & \dots & d_{m/2+1}^{i,j} \\ \vdots & \vdots & \vdots & \vdots \\ d_{m/2}^{i,j} & \dots & \dots & d_m^{i,j} \end{bmatrix}$$



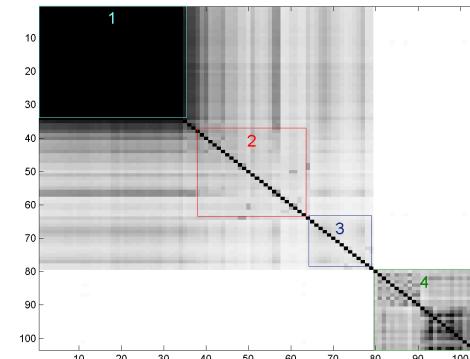
Hankel



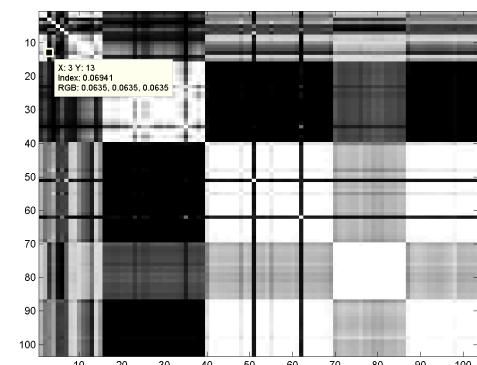
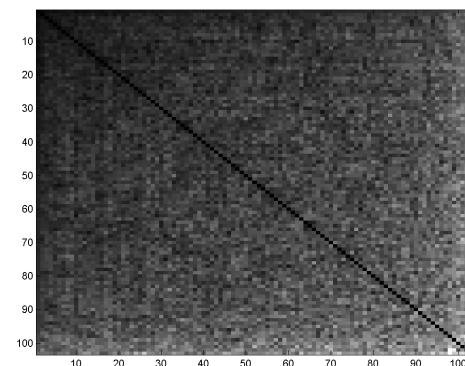
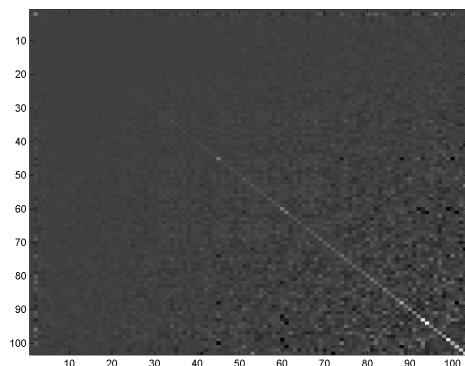
Costeira Kanade

Irani

GPCA



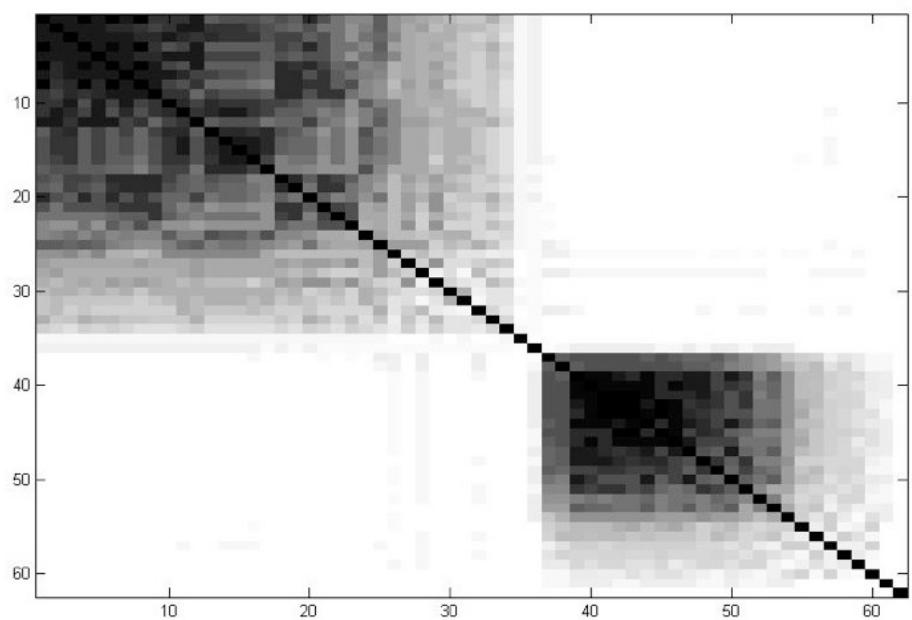
Hankel



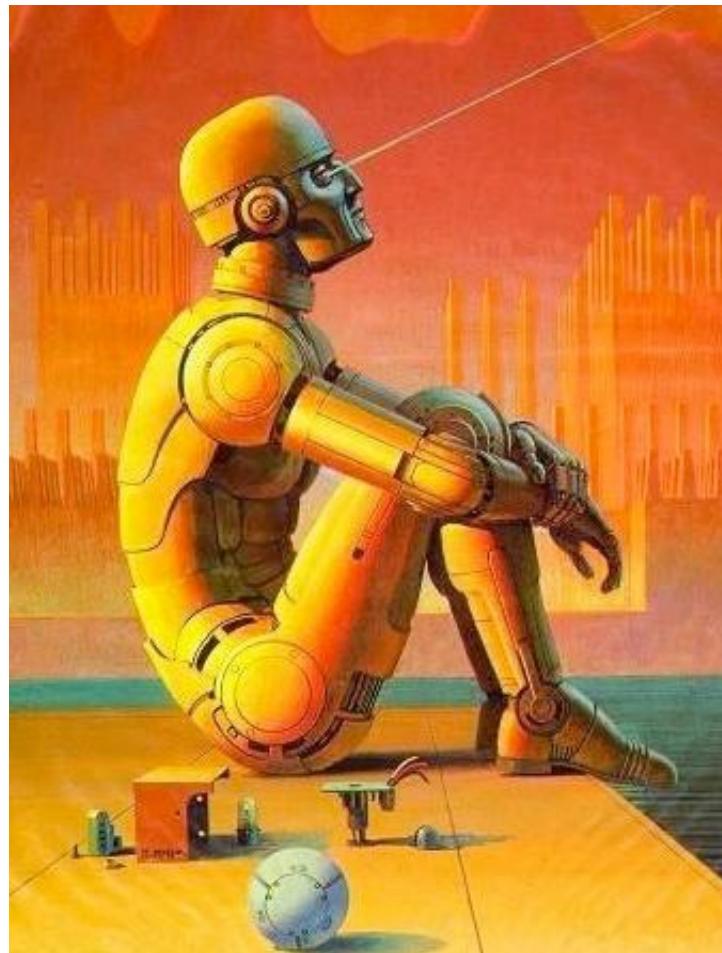
Costeira Kanade

Irani

GPCA



# 3D Object Recognition



Credit: some slides have been adopted from Fergus, Torralba, Fei, Fidler, and Grauman

# The Ultimate Goal



# Applications:

---

Industrial inspection, quality control

Surveillance and security

Assisted living

Human-computer interfaces

Medical image analysis

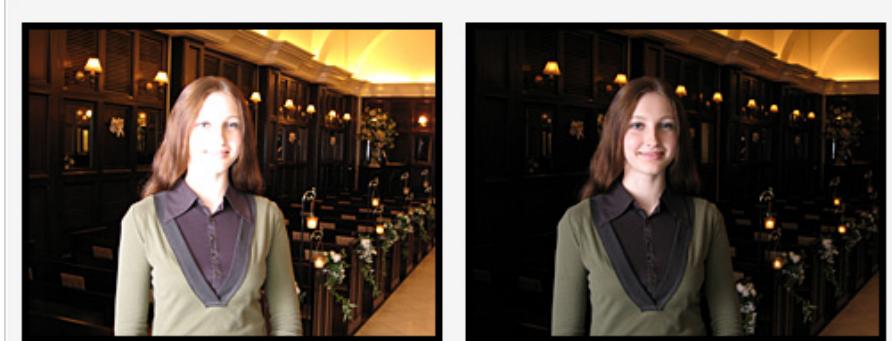
Reverse engineering

Image databases

# More Applications:



## Computational Photography

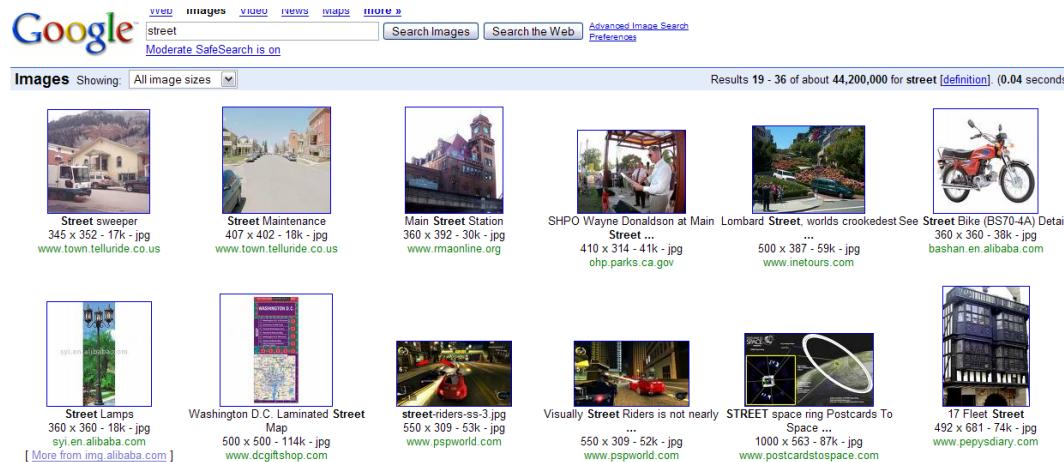


[Face priority AE] When a bright part of the face is too bright

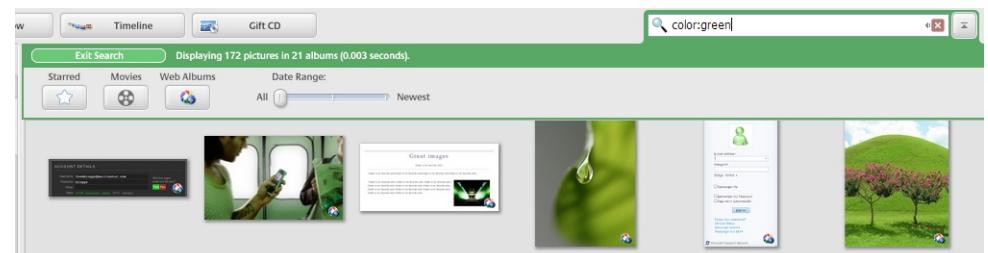
# Improving online search



Query:  
STREET

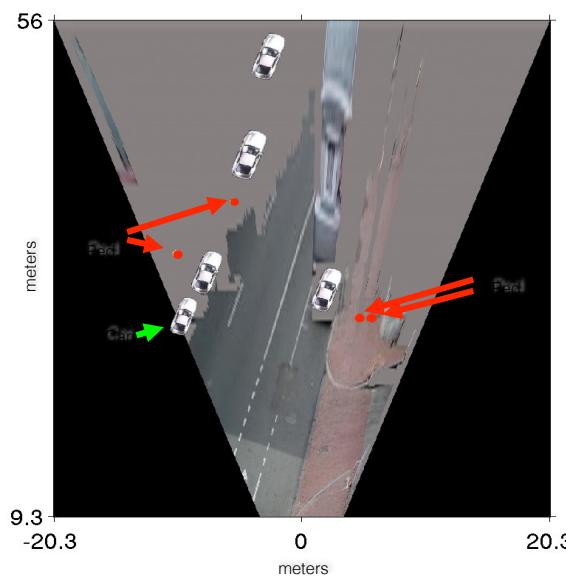


Organizing photo collections

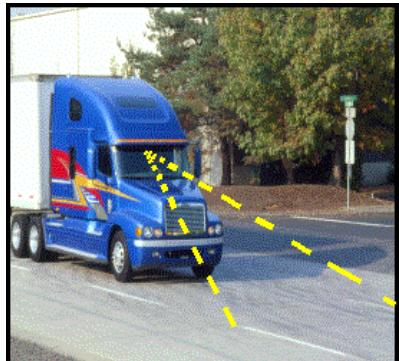


# Assisted driving

## Pedestrian and car detection



## Lane detection



- Collision warning systems with adaptive cruise control
- Lane departure warning systems
- Rear object detection systems

# What does Object Recognition Involve?

---



# What does Object Recognition Involve?

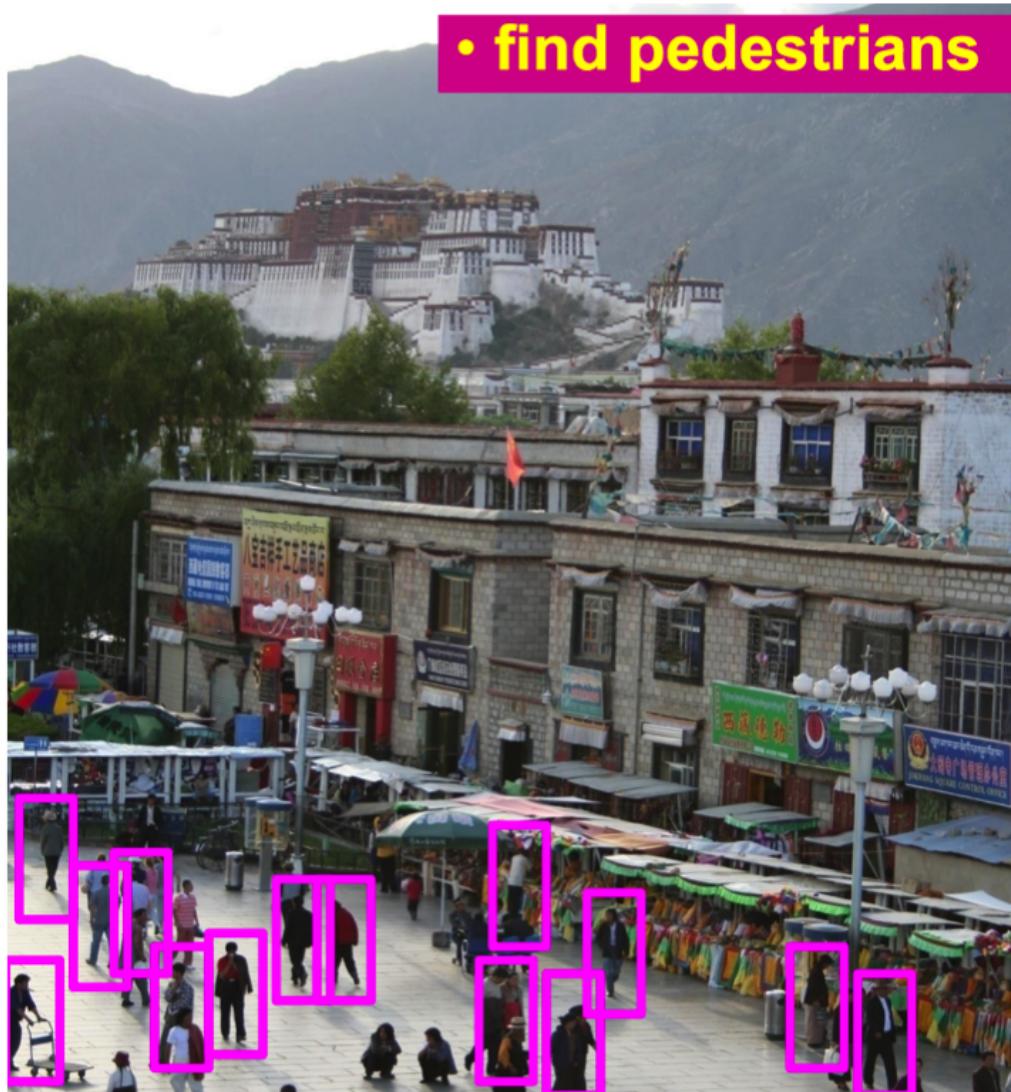
Detection:  
Are there people?



# What does Object Recognition Involve?

Detection:

Where are the people?



# What does Object Recognition Involve?



Verification:  
Is this a lamp?

# What does Object Recognition Involve?



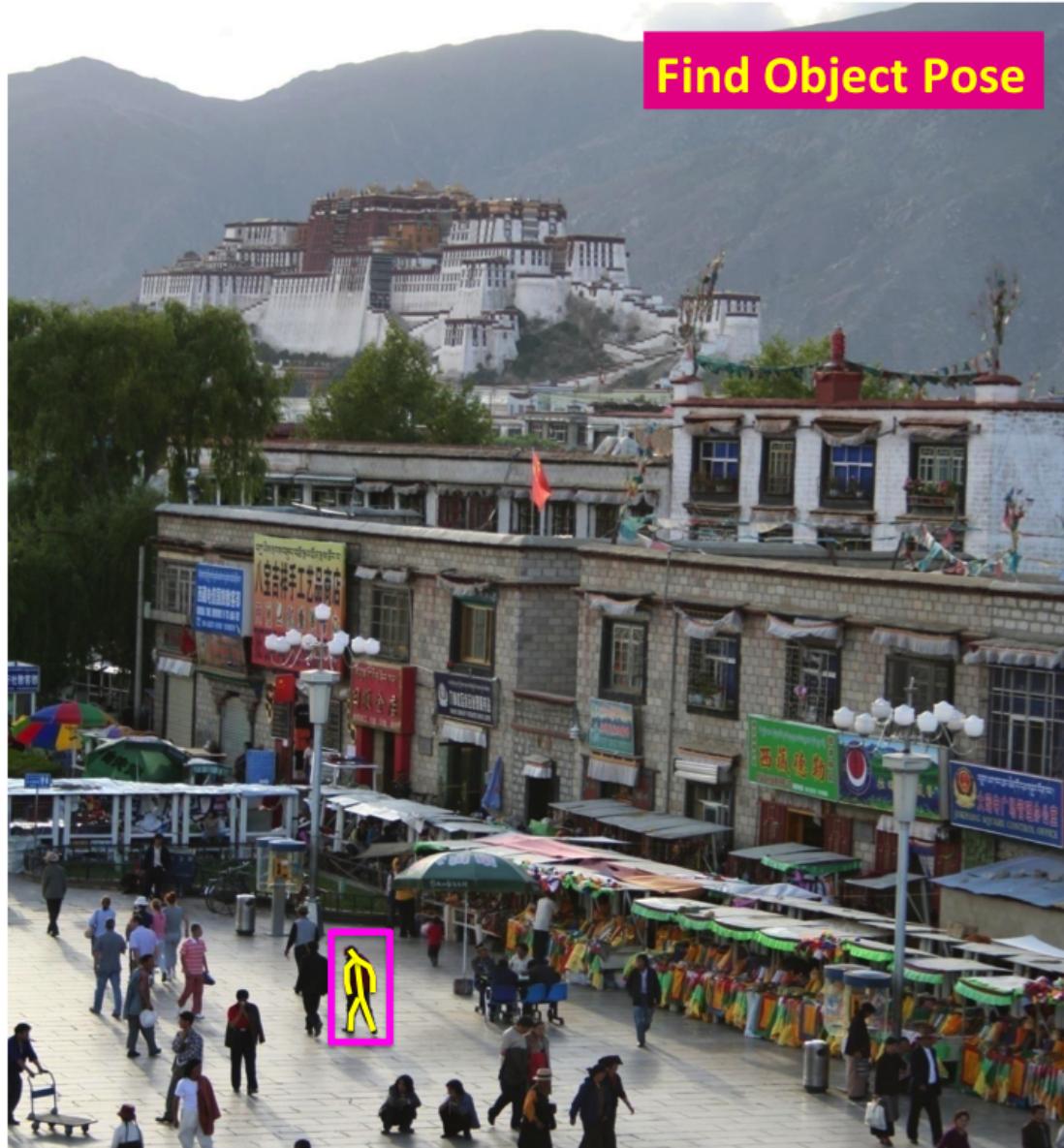
Identification:  
Is this Potala  
Palace?

# What does Object Recognition Involve?

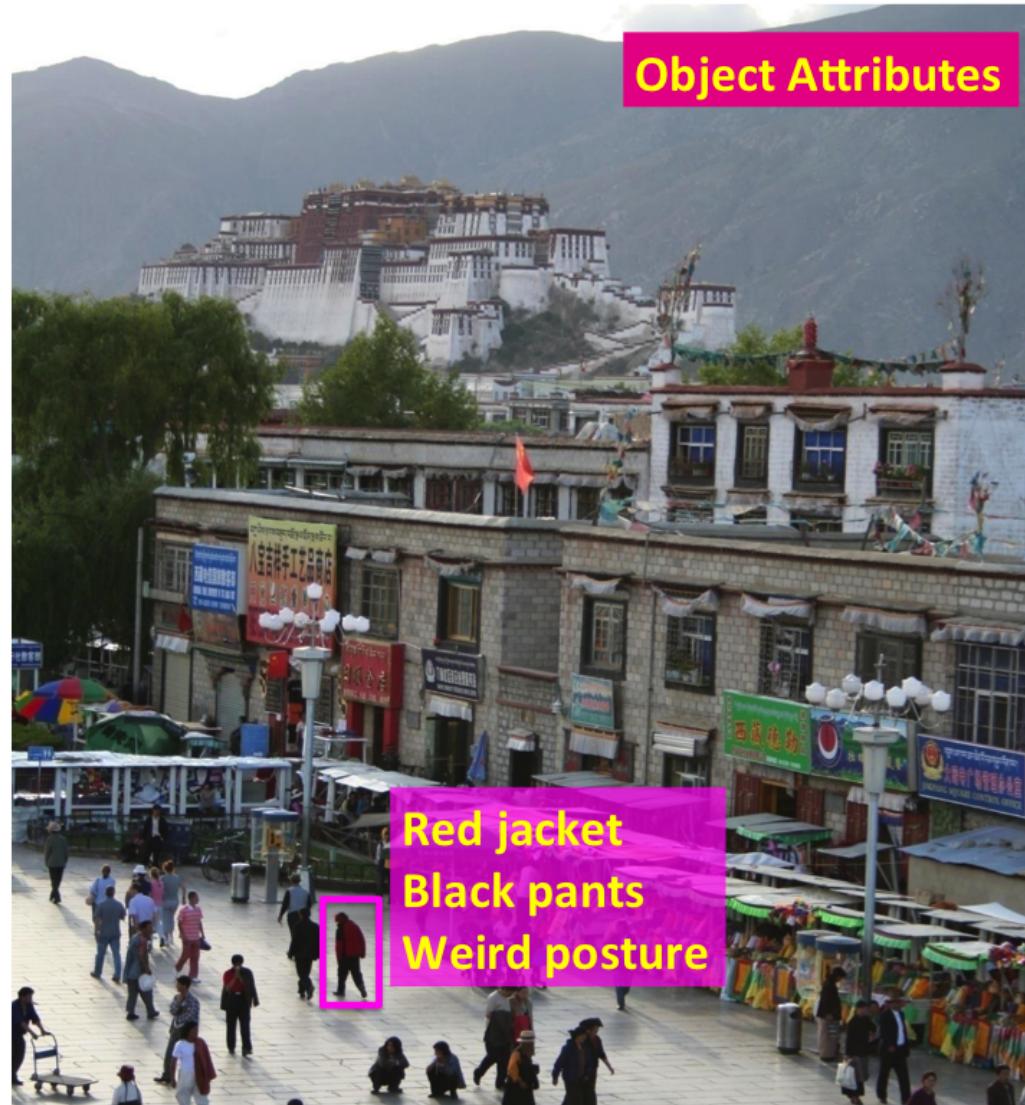


Object  
Categorization

# What does Object Recognition Involve?



# What does Object Recognition Involve?

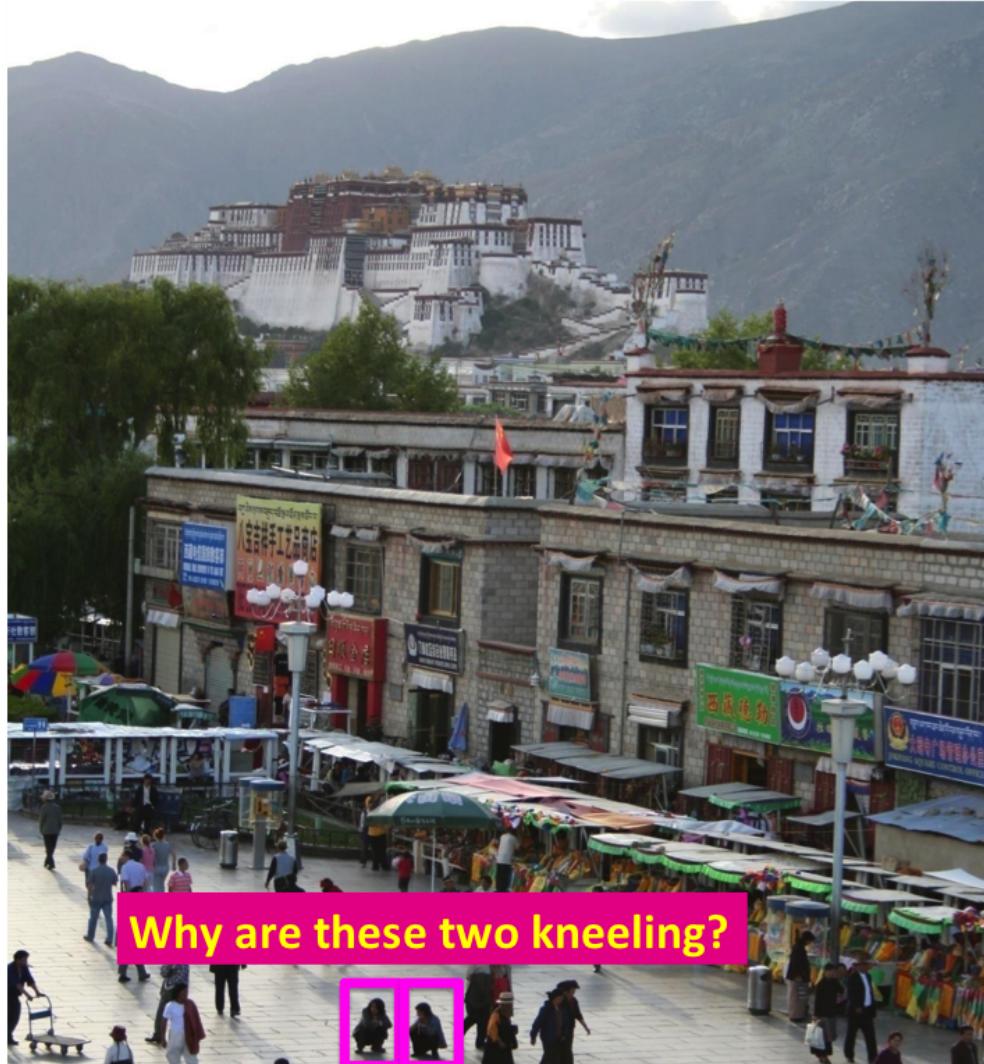


# What does Object Recognition Involve?

---



# What does Object Recognition Involve?



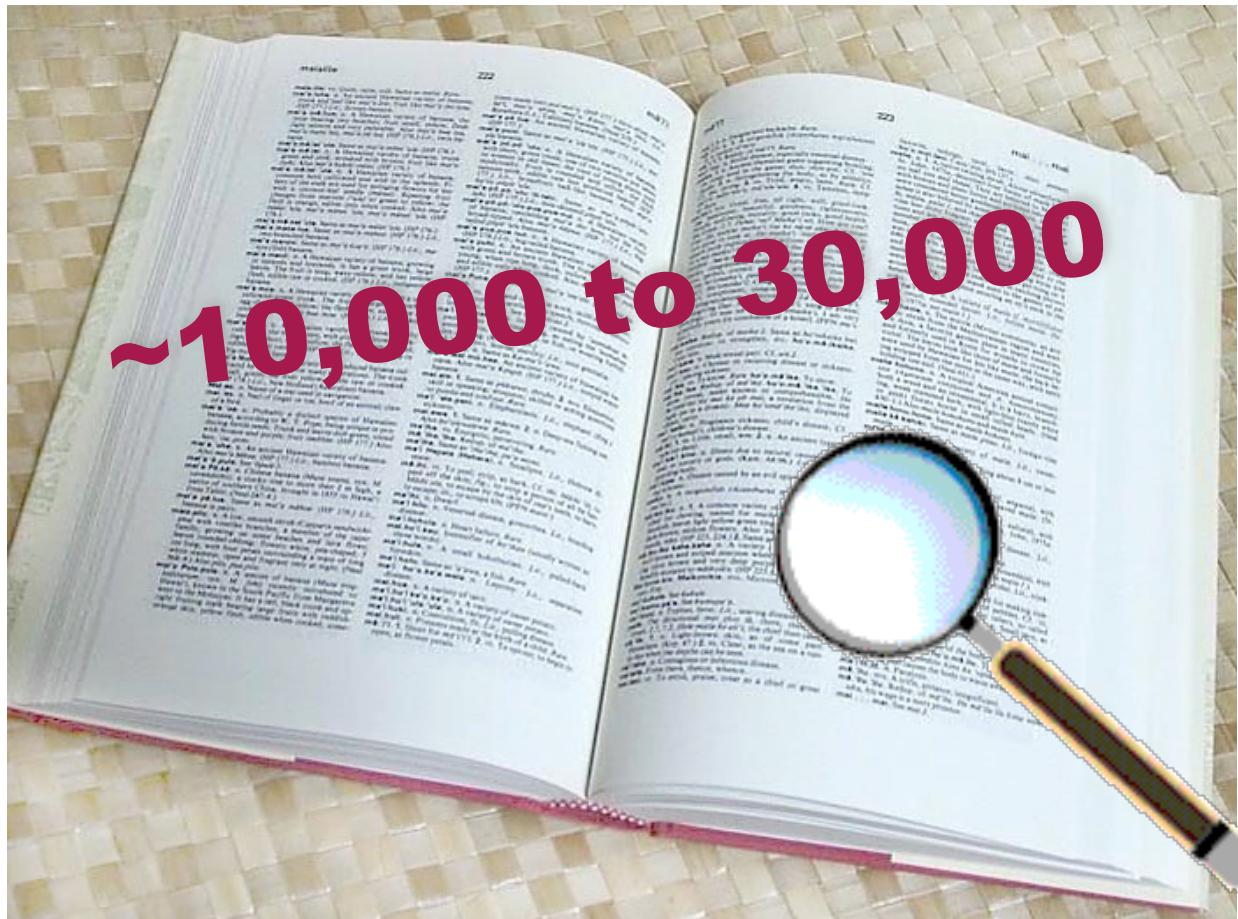
# What does Object Recognition Involve?



- **outdoor**
- **city**
- ...

Scene and  
Context  
Categorization

# How many object categories are there?



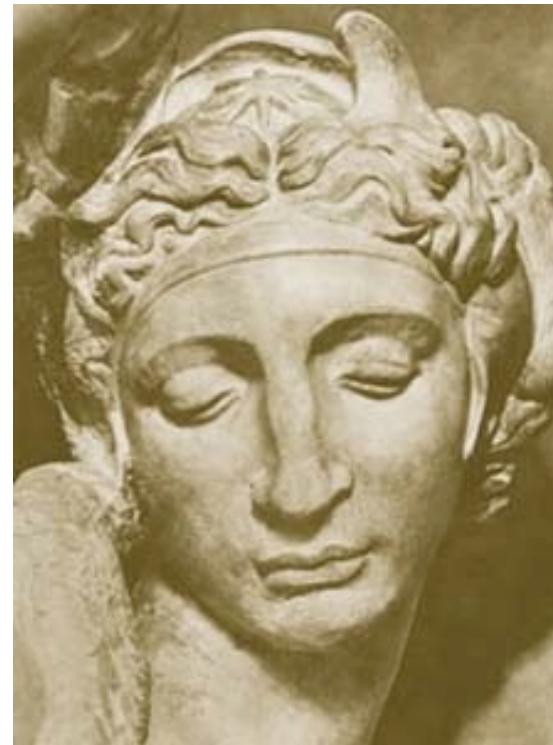
Biederman 1987

---

Not an easy task ...

# Challenges 1: view point variation

---



Michelangelo 1475-1564

## Challenges 2: illumination

---



# Challenges 3: occlusion

---



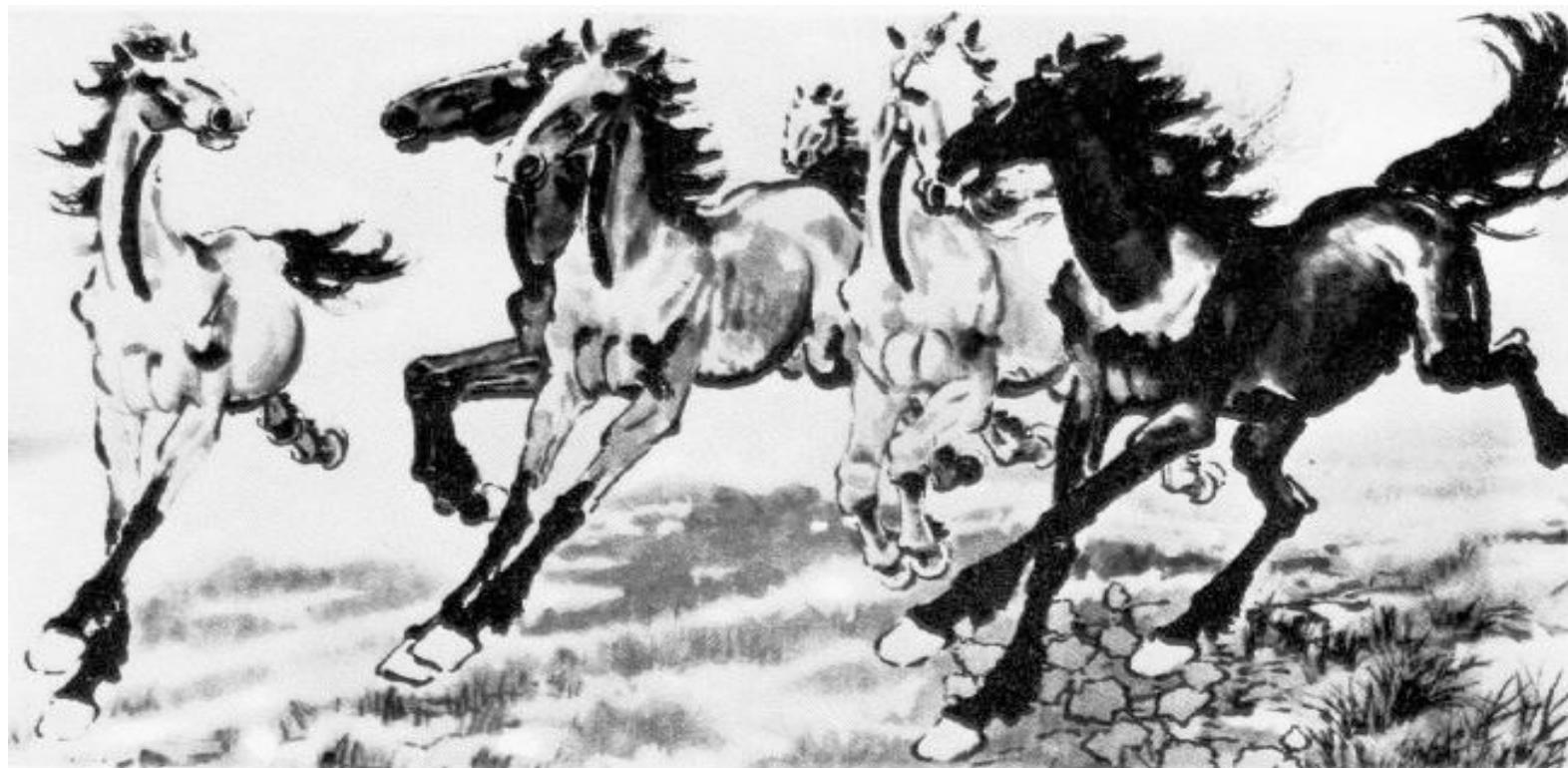
Magritte, 1957

# Challenges 4: scale



# Challenges 5: deformation

---



Xu, Beihong 1943

# Challenges 6: background clutter

---

Klimt, 1913

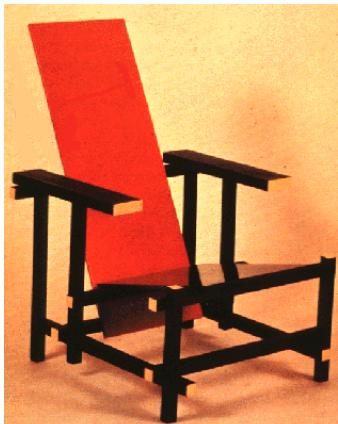


# History: single object recognition



# Challenges 7: intra-class variation

---





~10,000 to 30,000



# History: single object recognition

---



- Dickinson, 1991
- Camps and Shapiro, 1992
- Sclaroff, 1993
- Nayar, 1995
- Lowe, et al. 1999, 2003
- Mahamud and Herbert, 2000
- Ferrari, Tuytelaars, and Van Gool, 2004
- Rothganger, Lazebnik, and Ponce, 2004
- Moreels and Perona, 2005
- ...

# History: early object categorization



3	6	8	1	7	9	6	6	9	1
6	7	5	7	8	6	3	4	8	5
2	1	7	9	7	1	2	8	4	6
4	8	1	9	0	1	8	8	9	4
7	6	1	8	6	4	1	5	6	0
7	5	9	2	6	5	8	1	9	7
1	2	2	2	2	3	4	4	8	0
0	2	3	8	0	7	3	8	5	7
0	1	4	6	4	6	0	2	4	3
7	1	2	8	1	6	9	8	6	1



Turk and Pentland, 1991  
Belhumeur, Hespanha, & Kriegman, 1997  
Schneiderman & Kanade 2004  
Viola and Jones, 2000

4 8 1 9 0 1 8 8 9 4  
7 6 1 8 6 4 1 5 6 0  
7 5 9 2 6 5 8 1 9 7  
1 2 2 2 2 3 4 4 8 0  
0 2 3 8 0 7 3 8 5 7  
0 1 4 6 4 6 0 2 4 3  
7 1 2 8 1 6 9 8 6 1



- Amit and Geman, 1999
- LeCun et al. 1998
- Belongie and Malik, 2002
- Schneiderman & Kanade, 2004
- Argawal and Roth, 2002
- Poggio et al. 1993

# How would we approach this problem?

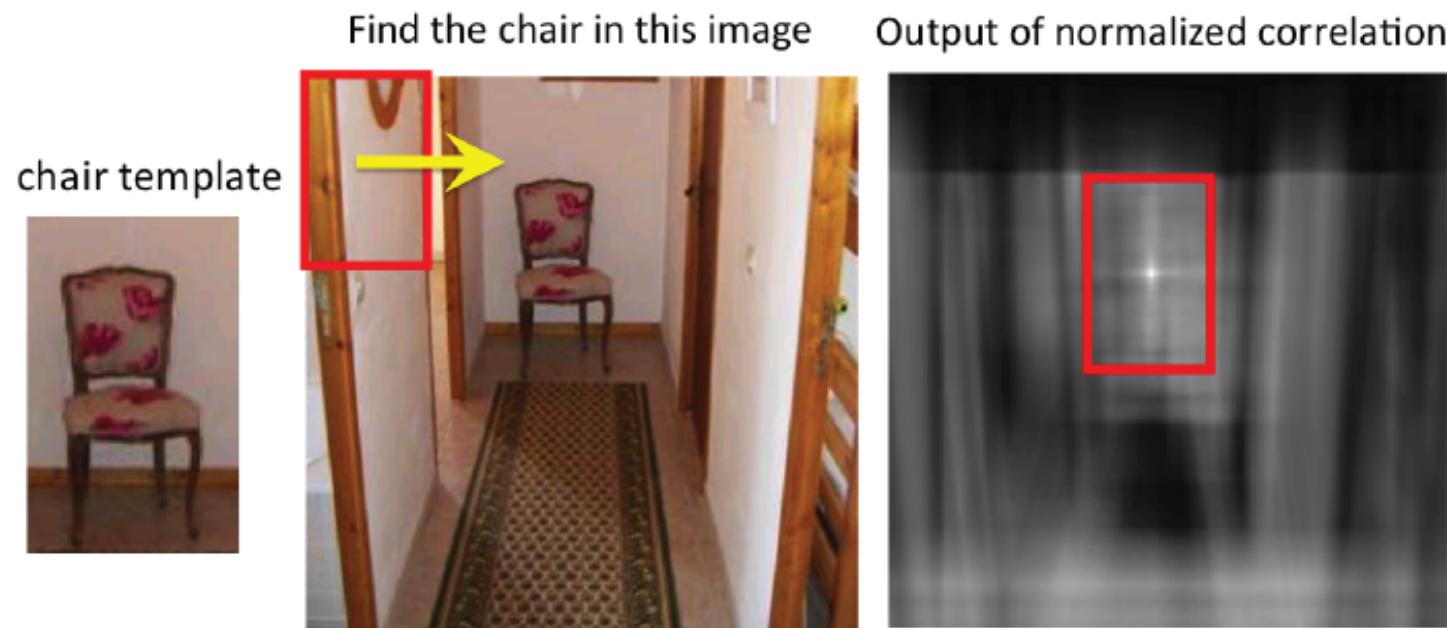
---

Let's try using what we already know to do DETECTION

Template Matching

# Template Matching

Normalize cross-correlation with a template (filter)



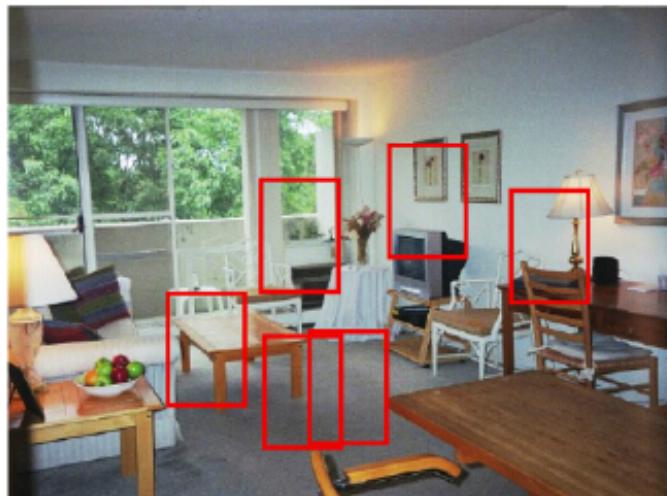
# Template Matching

Normalize cross-correlation with a template (filter)

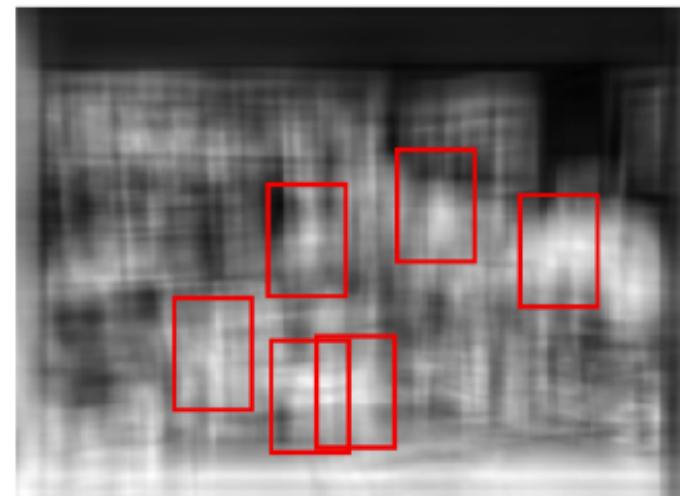


template

Find the chair in this image



Pretty much garbage  
Simple template matching is  
not going to make it



# How would we approach this problem?

---

Let's try using what we already know to do DETECTION

Large sale Retrieval

Store lots of images, recognize one by finding the most similar in the database. Google Approach.

# Recognition via Retrieval



query



# Recognition via Retrieval

---

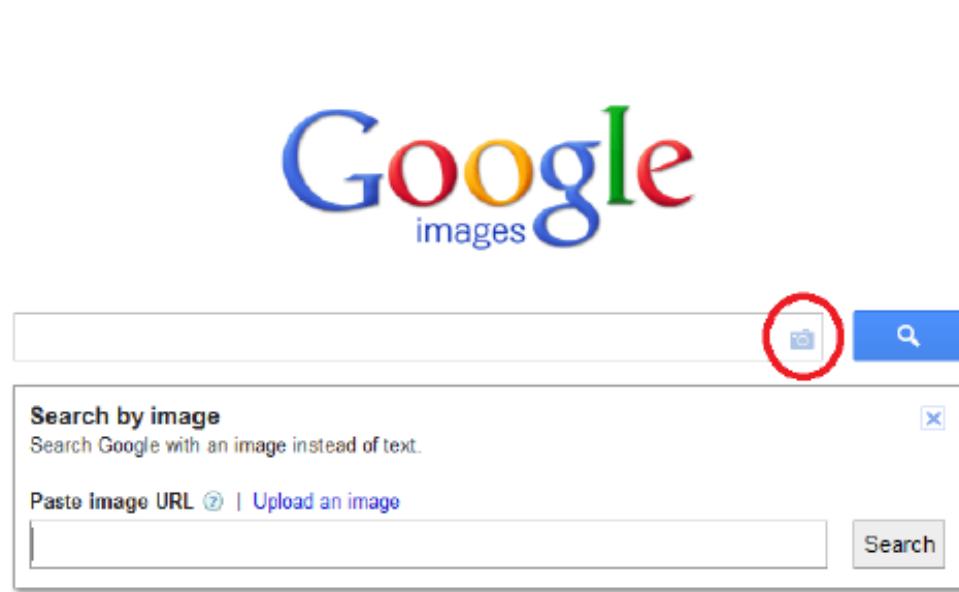


query



reasonable result!

# Recognition via Retrieval



query



# Recognition via Retrieval

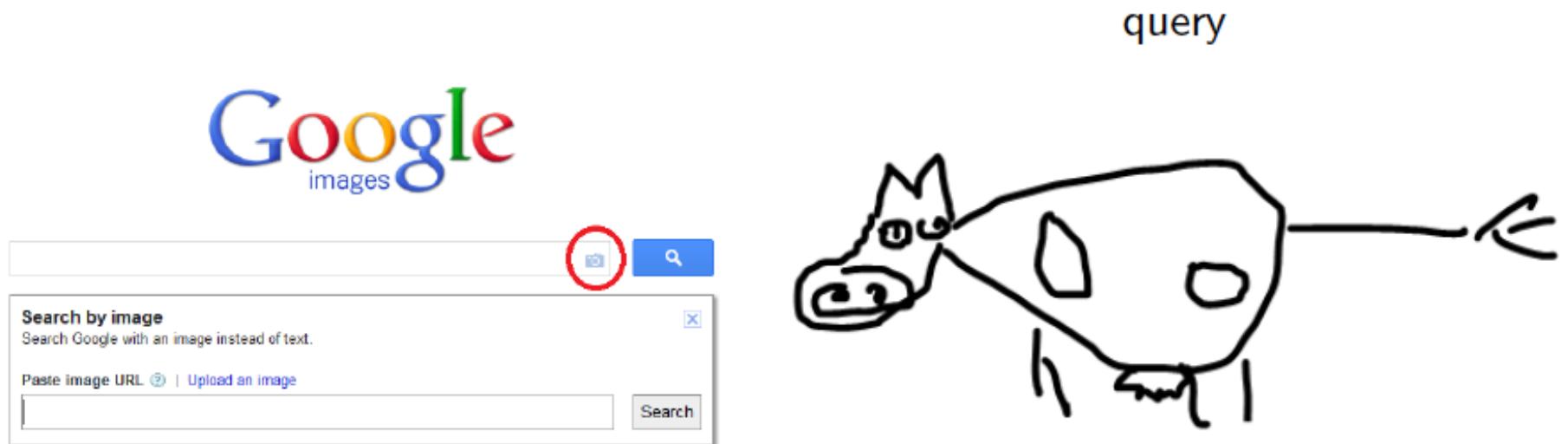


query

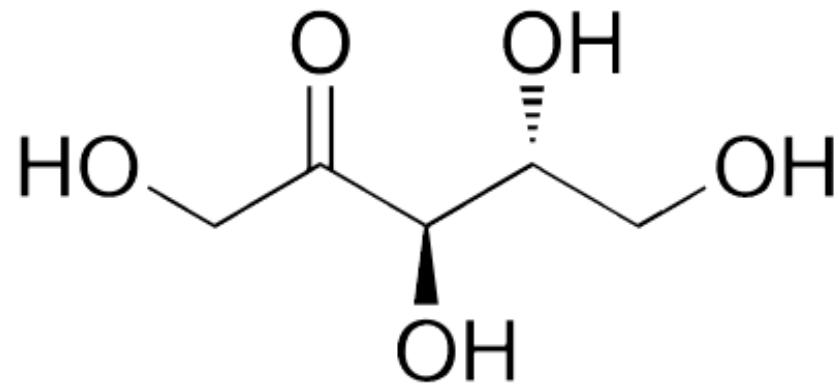


less reasonable result, but somewhat similar?

# Recognition via Retrieval



# Recognition via Retrieval



not a reasonable result!

# Indexing for Fast Retrieval

---

PCA: Eigenimages

Local Features: Bag of Words

# Template Matching

---

Objects can be represented by storing sample images or “templates”

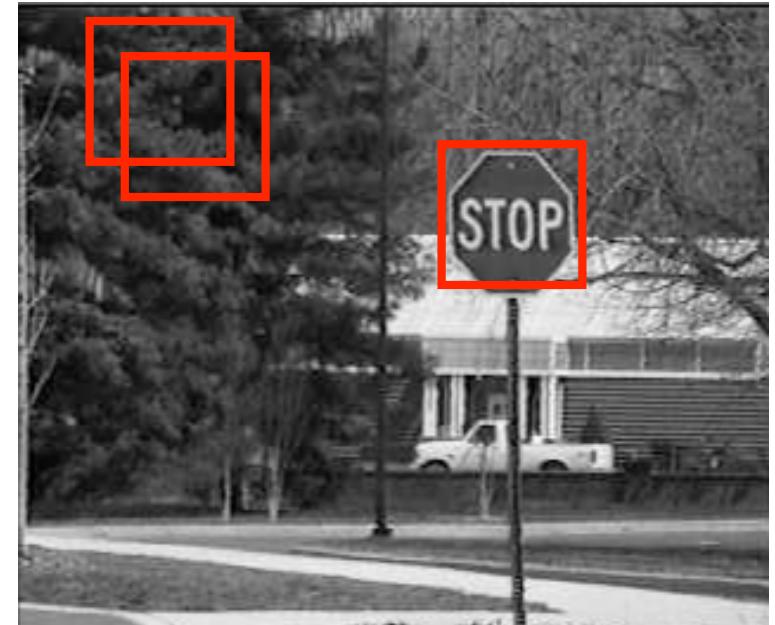


Stop sign template

# Hypotheses from Template Matching



- Place the template at every location on the given image.
- Compare the pixel values in the template with the pixel values in the underlying region of the image.
- If a “good” match is found, announce that the object is present in the image.



- Possible measures are: SSD, SAD, Cross-correlation, Normalized Cross-correlation, max difference, etc.

# Limitations of Template Matching

If the object appears scaled, rotated, or skewed on the image, the match will not be good.



# Solution:

---

Search for the template and possible transformations of the template:



Not very efficient! (but doable ...)

# Limitations of Template Matching

---

It uses global information: it is sensitive to occlusion.



# Limitations of Template Matching

---

It uses pixel values: it is illumination and sensor dependent.



# Eigenimages

---

The appearance of an object in an image depends on several things:

Viewpoint

Illumination conditions

Sensor

The object itself (ex: human facial expression)

In principle, these variations can be handled by increasing the number of templates.

# Eigenimages: Using multiple templates

---

- The number of templates can grow very fast!
- We need:
  - An efficient way to store templates
  - An efficient way to search for matches
- Observation: while each template is different, there exist many similarities between the templates.

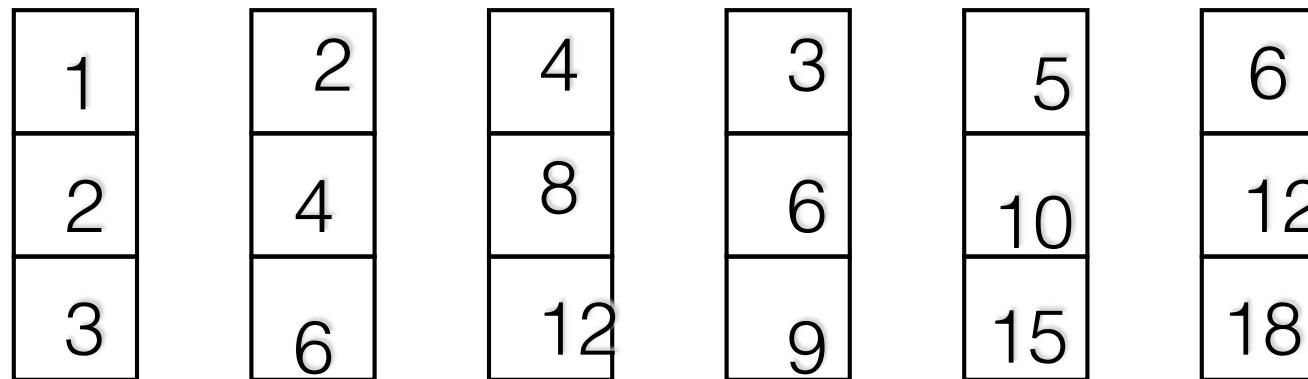


# Efficient Image Storage

---

Toy Example: Images with 3 pixels

Consider the following 3x1 templates:



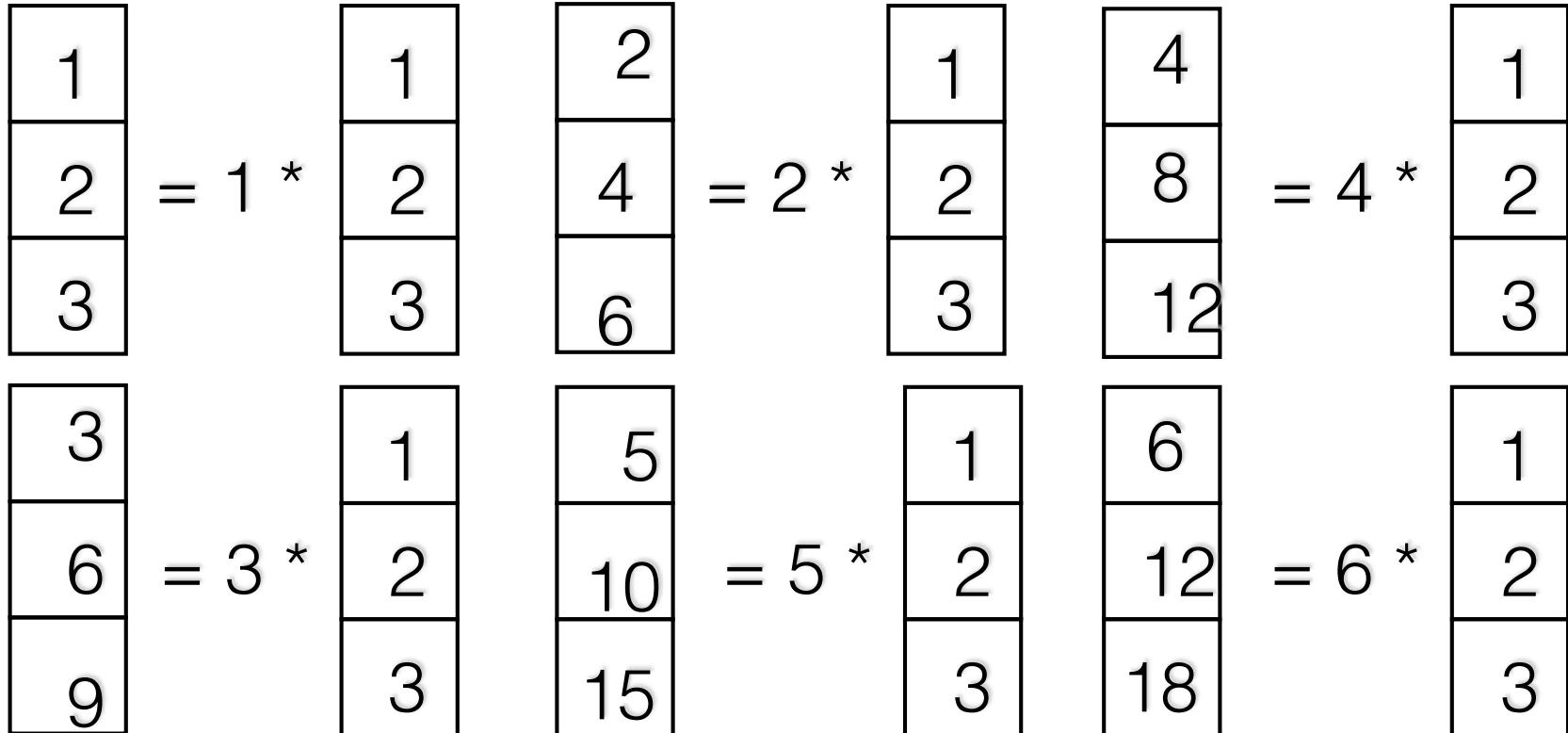
If each pixel is stored in a byte, we need  $3 \times 6 = 18$  bytes

# Efficient Image Storage

Looking closer, we can see that all the images are very similar to each other: they are all the same image, scaled by a factor:

$$\begin{array}{c|c} \begin{array}{|c|c|c|} \hline 1 & & \\ \hline 2 & & \\ \hline 3 & & \\ \hline \end{array} & = 1 * \begin{array}{|c|c|} \hline 1 & \\ \hline 2 & \\ \hline 3 & \\ \hline \end{array} \end{array}$$
$$\begin{array}{c|c} \begin{array}{|c|c|c|} \hline 2 & & \\ \hline 4 & & \\ \hline 6 & & \\ \hline \end{array} & = 2 * \begin{array}{|c|c|} \hline 1 & \\ \hline 2 & \\ \hline 3 & \\ \hline \end{array} \end{array}$$
$$\begin{array}{c|c} \begin{array}{|c|c|c|} \hline 4 & & \\ \hline 8 & & \\ \hline 12 & & \\ \hline \end{array} & = 4 * \begin{array}{|c|c|} \hline 1 & \\ \hline 2 & \\ \hline 3 & \\ \hline \end{array} \end{array}$$
$$\begin{array}{c|c} \begin{array}{|c|c|c|} \hline 1 & & \\ \hline 2 & & \\ \hline 3 & & \\ \hline \end{array} & \end{array}$$
$$\begin{array}{c|c} \begin{array}{|c|c|c|} \hline 3 & & \\ \hline 6 & & \\ \hline 9 & & \\ \hline \end{array} & = 3 * \begin{array}{|c|c|} \hline 1 & \\ \hline 2 & \\ \hline 3 & \\ \hline \end{array} \end{array}$$
$$\begin{array}{c|c} \begin{array}{|c|c|c|} \hline 5 & & \\ \hline 10 & & \\ \hline 15 & & \\ \hline \end{array} & = 5 * \begin{array}{|c|c|} \hline 1 & \\ \hline 2 & \\ \hline 3 & \\ \hline \end{array} \end{array}$$
$$\begin{array}{c|c} \begin{array}{|c|c|c|} \hline 1 & & \\ \hline 2 & & \\ \hline 3 & & \\ \hline \end{array} & \end{array}$$
$$\begin{array}{c|c} \begin{array}{|c|c|c|} \hline 6 & & \\ \hline 12 & & \\ \hline 18 & & \\ \hline \end{array} & = 6 * \begin{array}{|c|c|} \hline 1 & \\ \hline 2 & \\ \hline 3 & \\ \hline \end{array} \end{array}$$
$$\begin{array}{c|c} \begin{array}{|c|c|c|} \hline 1 & & \\ \hline 2 & & \\ \hline 3 & & \\ \hline \end{array} & \end{array}$$

# Efficient Image Storage

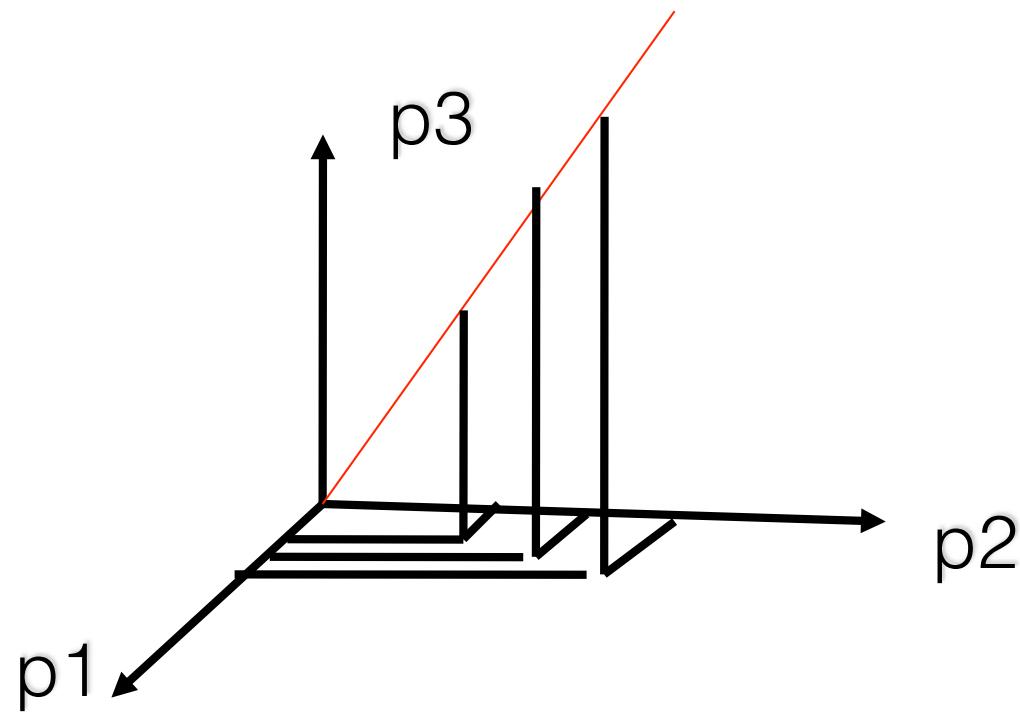


They can be stored using only 9 bytes (50% savings!):  
Store one image (3 bytes) + the multiplying constants (6 bytes)

# Geometrical Interpretation:

---

Consider each pixel in the image as a coordinate in a vector space. Then, each  $3 \times 1$  template can be thought of as a point in a 3D space:

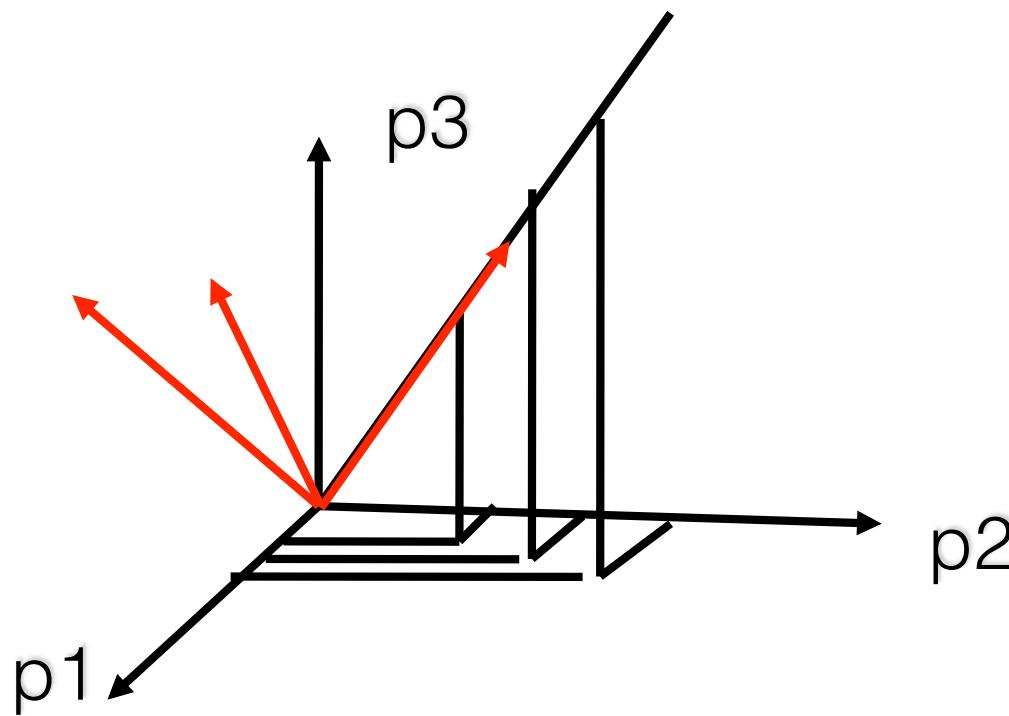


But in this example, all the points happen to belong to a line: a 1D subspace of the original 3D space.

# Geometrical Interpretation:

---

Consider a new coordinate system where one of the axes is along the direction of the line:



In this coordinate system, every image has only one non-zero coordinate: we only need to store the direction of the line (a 3 bytes image) and the non-zero coordinate for each of the images (6 bytes).

# Principal Component Analysis (PCA)

---

Given a set of templates, how do we know if they can be compressed like in the previous example?

The answer is to look into the correlation between the templates

The tool for doing this is called PCA

# PCA Theorem

---

Let  $x_1 x_2 \dots x_n$  be a set of  $n N^2 \times 1$  vectors and let  $\bar{x}$  be their average:

$$x_i = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{iN^2} \end{bmatrix} \quad \bar{x} = \frac{1}{n} \sum_{i=1}^{i=n} \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{iN^2} \end{bmatrix}$$

**Note:** Each  $N \times N$  image template can be represented as a  $N^2 \times 1$  vector whose elements are the template pixel values.

# PCA Theorem

---

Let  $X$  be the  $N^2 \times n$  matrix:

$$X = \begin{bmatrix} \mathbf{x}_1 - \bar{\mathbf{x}} & \mathbf{x}_2 - \bar{\mathbf{x}} & \cdots & \mathbf{x}_n - \bar{\mathbf{x}} \end{bmatrix}$$

**Note:** subtracting the mean is equivalent to translating the coordinate system to the location of the mean.

# PCA Theorem

---

Let  $Q = X X^T$  be the  $N^2 \times N^2$  matrix:

$$Q = X X^T = \begin{bmatrix} \mathbf{x}_1 - \bar{\mathbf{x}} & \mathbf{x}_2 - \bar{\mathbf{x}} & \cdots & \mathbf{x}_n - \bar{\mathbf{x}} \end{bmatrix} \begin{bmatrix} (\mathbf{x}_1 - \bar{\mathbf{x}})^T \\ (\mathbf{x}_2 - \bar{\mathbf{x}})^T \\ \vdots \\ (\mathbf{x}_n - \bar{\mathbf{x}})^T \end{bmatrix}$$

## Notes:

1.  $Q$  is square
2.  $Q$  is symmetric
3.  $Q$  is the covariance matrix
4.  $Q$  can be very large (remember that  $N^2$  is the number of pixels in the template)

# PCA Theorem

---

Theorem:

Each  $x_j$  can be written as:  $x_j = \bar{x} + \sum_{i=1}^{i=n} g_{ji} e_i$

where  $e_i$  are the n eigenvectors of Q with non-zero eigenvalues.

## Notes:

1. The eigenvectors  $e_1 e_2 \dots e_n$  span an **eigenspace**
2.  $e_1 e_2 \dots e_n$  are  $N^2 \times 1$  orthonormal vectors ( $N \times N$  images).
3. The scalars  $g_{ji}$  are the coordinates of  $x_j$  in the space.
4.  $g_{ji} = (x_j - \bar{x}) \cdot e_i$

# Using PCA to Compress Data

---

Expressing  $x$  in terms of  $e_1 \dots e_n$  has not changed the size of the data

However, if the templates are highly correlated many of the coordinates of  $x$  will be zero or close to zero.

# Using PCA to Compress Data

---

Sort the eigenvectors  $e_i$  according to their eigenvalue:

$$\lambda_1 \geq \lambda_2 \geq \dots \lambda_n$$

- Assuming that  $\lambda_i \approx 0$  if  $i > k$

- Then

$$\mathbf{x}_j \approx \bar{\mathbf{x}} + \sum_{i=1}^{i=k} g_{ji} \mathbf{e}_i$$

# Eigenspaces: Efficient Image Storage



- Use PCA to compress the data:
  - each image is stored as a k-dimensional vector
  - Need to store  $k N \times N$  eigenvectors
  - $k \ll n \ll N^2$

$$\begin{matrix} \text{[Truck Image]} & \cong & \text{[Truck Image]} & = a_{01} & \text{[Eigenvector]} & + a_{02} & \text{[Eigenvector]} & + a_{03} & \text{[Eigenvector]} & + a_{04} & \text{[Eigenvector]} & + a_{05} & \text{[Eigenvector]} & + a_{06} & \text{[Eigenvector]} & + \dots \end{matrix}$$

# Eigenspaces: Efficient Image Comparison



- Use the same procedure to compress the given image to a  $k$ -dimensional vector.
- Compare the compressed vectors:
  - Dot product of  $k$ -dimensional vectors
  - $k \ll n \ll N^2$

$$\begin{matrix} \text{[Truck Image]} & \cong & \text{[Truck Image]} & = a_{01} & \text{[Eigenvector 1]} & + a_{02} & \text{[Eigenvector 2]} & + a_{03} & \text{[Eigenvector 3]} & + a_{04} & \text{[Eigenvector 4]} & + a_{05} & \text{[Eigenvector 5]} & + a_{06} & \text{[Eigenvector 6]} & + \dots \end{matrix}$$

# Implementing PCA

---

Need to find “first” k eigenvectors of Q:

$$Q = XX^T = \begin{bmatrix} \mathbf{x}_1 - \bar{\mathbf{x}} & \mathbf{x}_2 - \bar{\mathbf{x}} & \cdots & \mathbf{x}_n - \bar{\mathbf{x}} \end{bmatrix} \begin{bmatrix} (\mathbf{x}_1 - \bar{\mathbf{x}})^T \\ (\mathbf{x}_2 - \bar{\mathbf{x}})^T \\ \vdots \\ (\mathbf{x}_n - \bar{\mathbf{x}})^T \end{bmatrix}$$

Q is  $N^2 \times N^2$  where  $N^2$  is the number of pixels in each image.  
For a 256 x 256 image,  $N^2 = 65536 !!$

# Finding ev of Q

---

$Q=XX^T$  is very large. Instead, consider the matrix  $P=X^TX$

- Q and P are both symmetric, but  $Q \neq P^T$
- Q is  $N^2 \times N^2$ , P is  $n \times n$
- n is the number of training images, typically  $n \ll N$

# Finding ev of Q

---

Let  $e$  be an eigenvector of  $P$  with eigenvalue  $\lambda$ :

$$Pe = \lambda e$$

$$X^T X e = \lambda e$$

$$XX^T X e = \lambda X e$$

$$Q(Xe) = \lambda(Xe)$$

$Xe$  is an eigenvector of  $Q$  also with eigenvalue  $\lambda$ !

# Singular Value Decomposition (SVD)

---

Any  $m \times n$  matrix  $X$  can be written as the product of 3 matrices:

$$X = UDV^T$$

Where:

- $U$  is  $m \times m$  and its columns are orthonormal vectors
- $V$  is  $n \times n$  and its columns are orthonormal vectors
- $D$  is  $m \times n$  diagonal and its diagonal elements are called the singular values of  $X$ , and are such that:

$$\sigma_1, \sigma_2, \dots, \sigma_n >= 0$$

# SVD Properties

---

$$X = UDV^T$$

- The columns of  $U$  are the eigenvectors of  $XX^T$
- The columns of  $V$  are the eigenvectors of  $X^TX$
- The squares of the diagonal elements of  $D$  are the eigenvalues of  $XX^T$  and  $X^TX$

# Algorithm EIGENSPACE\_LEARN

---

## Assumptions:

1. Each image contains one object only.
2. Objects are imaged by a fixed camera .
3. Images are normalized in size N x N:  
The image frame is the minimum rectangle enclosing the object.
4. Energy of pixels values is normalized to 1:  
 $\sum_i \sum_j I(i,j)^2 = 1$
5. The object is completely visible and unoccluded in all images.

# Algorithm EIGENSPACE\_LEARN

---

## Getting the data:

For each object  $o$  to be represented,  $o = 1, \dots, O$

1. Place  $o$  on a turntable, acquire a set of  $n$  images by rotating the table in increments of  $360^\circ/n$
2. For each image  $p$ ,  $p = 1, \dots, n$ :
  1. Segment  $o$  from the background
  2. Normalize the image size and energy
  3. Arrange the pixels as vectors  $\mathbf{x}_p^o$

# Algorithm EIGENSPACE\_LEARN

---

## Storing the data:

1. Find the average image vector  $\bar{\mathbf{x}} = \frac{1}{n.o} \sum_{o=1}^O \sum_{p=1}^n x_p^o$

2. Assemble the matrix X:

$$X = [ \mathbf{x}_1^1 - \bar{\mathbf{x}} \ \ \mathbf{x}_2^1 - \bar{\mathbf{x}} \ \ \dots \ \ \mathbf{x}_n^o - \bar{\mathbf{x}} ]$$

3. Find the first k eigenvectors of  $XX^T$ :  $e_1, \dots, e_k$   
(use  $X^TX$  or SVD)

4. For each object o, each image p:

• Compute the corresponding k-dimensional point:

$$\mathbf{g}_p^o = [ e_1 \ e_2 \ \dots \ e_k ] (\mathbf{x}_p^o - \bar{\mathbf{x}})$$

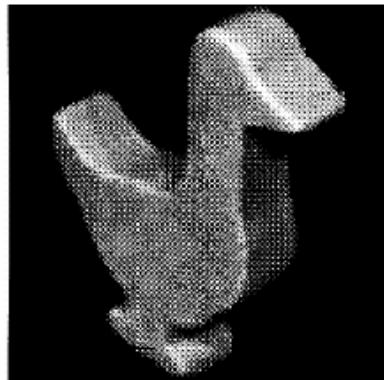
# Algorithm EIGENSPACE\_IDENTIF

---

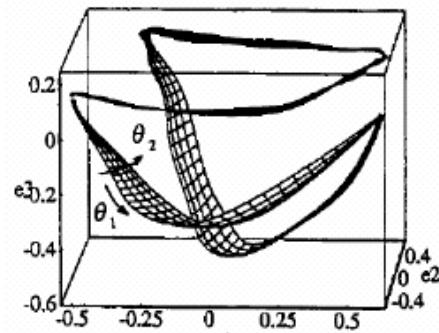
## Recognizing an object from the DB:

1. Given an image, segment the object from the background
2. Normalize the size an energy, write it as a vector  $\mathbf{i}$
3. Compute the corresponding k-dimensional point:
$$\mathbf{g} = [ \mathbf{e}_1 \quad \mathbf{e}_2 \quad \cdots \quad \mathbf{e}_k ] (\mathbf{i} - \bar{\mathbf{x}})$$
4. Find the closest  $\mathbf{g}^o_p$  k-dimensional point to  $\mathbf{g}$

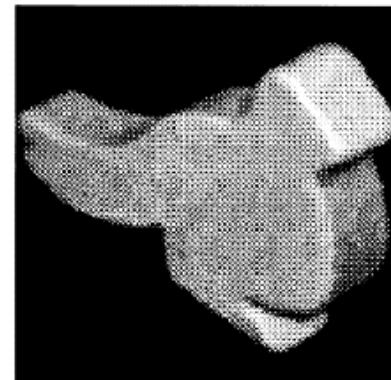
# Appearance Manifolds



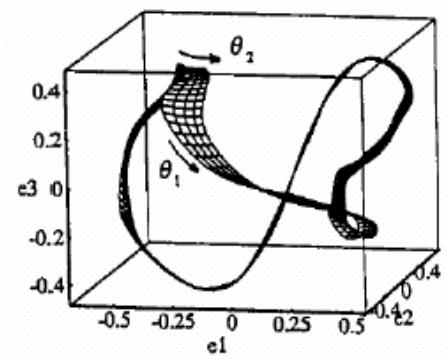
A



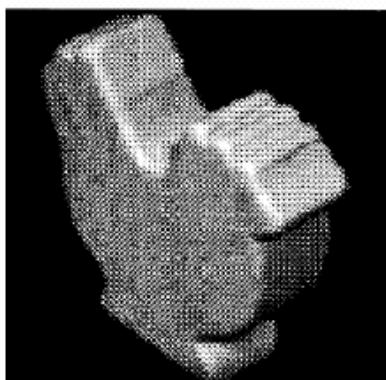
A



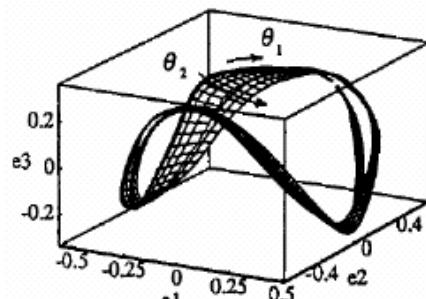
B



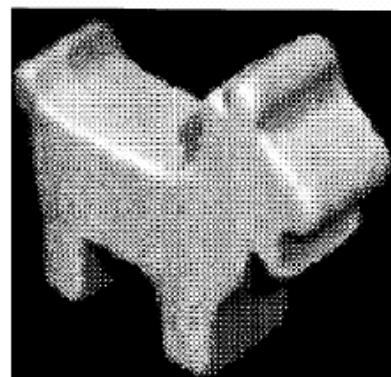
B



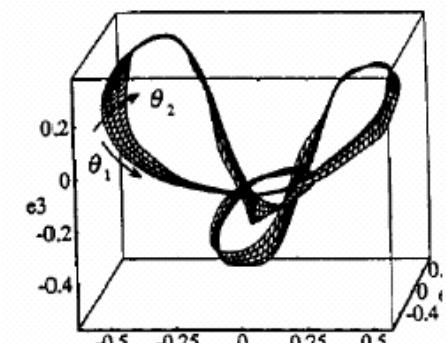
C



C



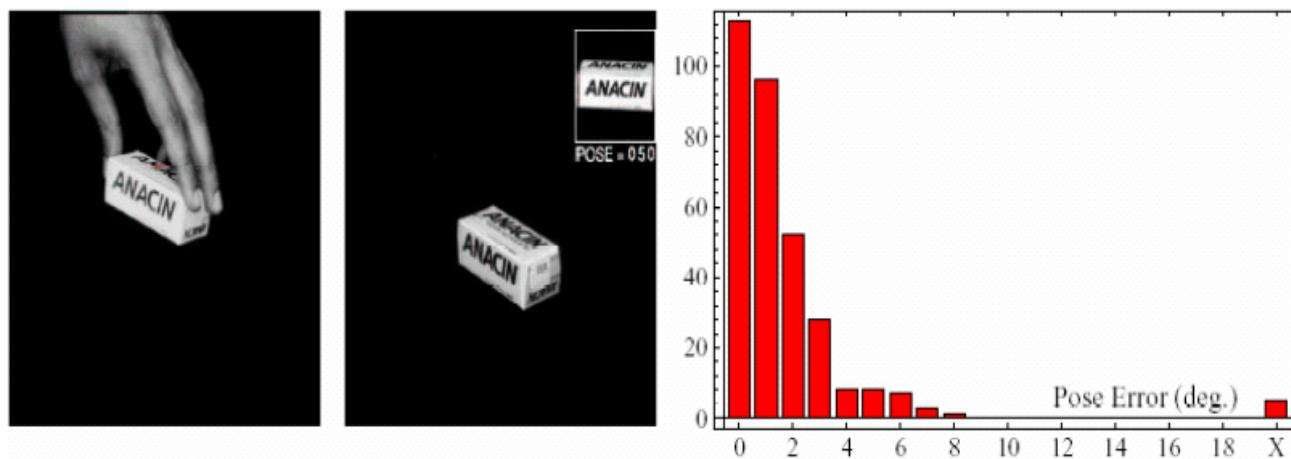
D



D

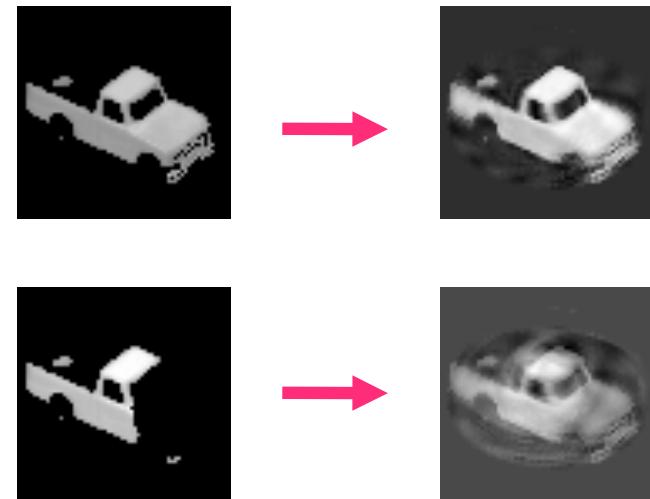
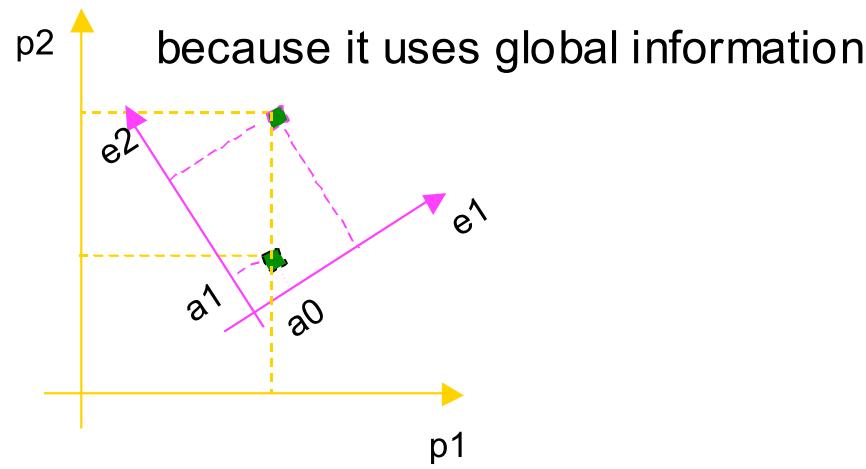
[Murase & Nayar, 1996]

# Appearance Manifolds



Murase & Nayar

# PCA has problems with occlusion



# Example 3: ABPs and ABRs

---

Use local appearance:

## PARTS and RELATIONS

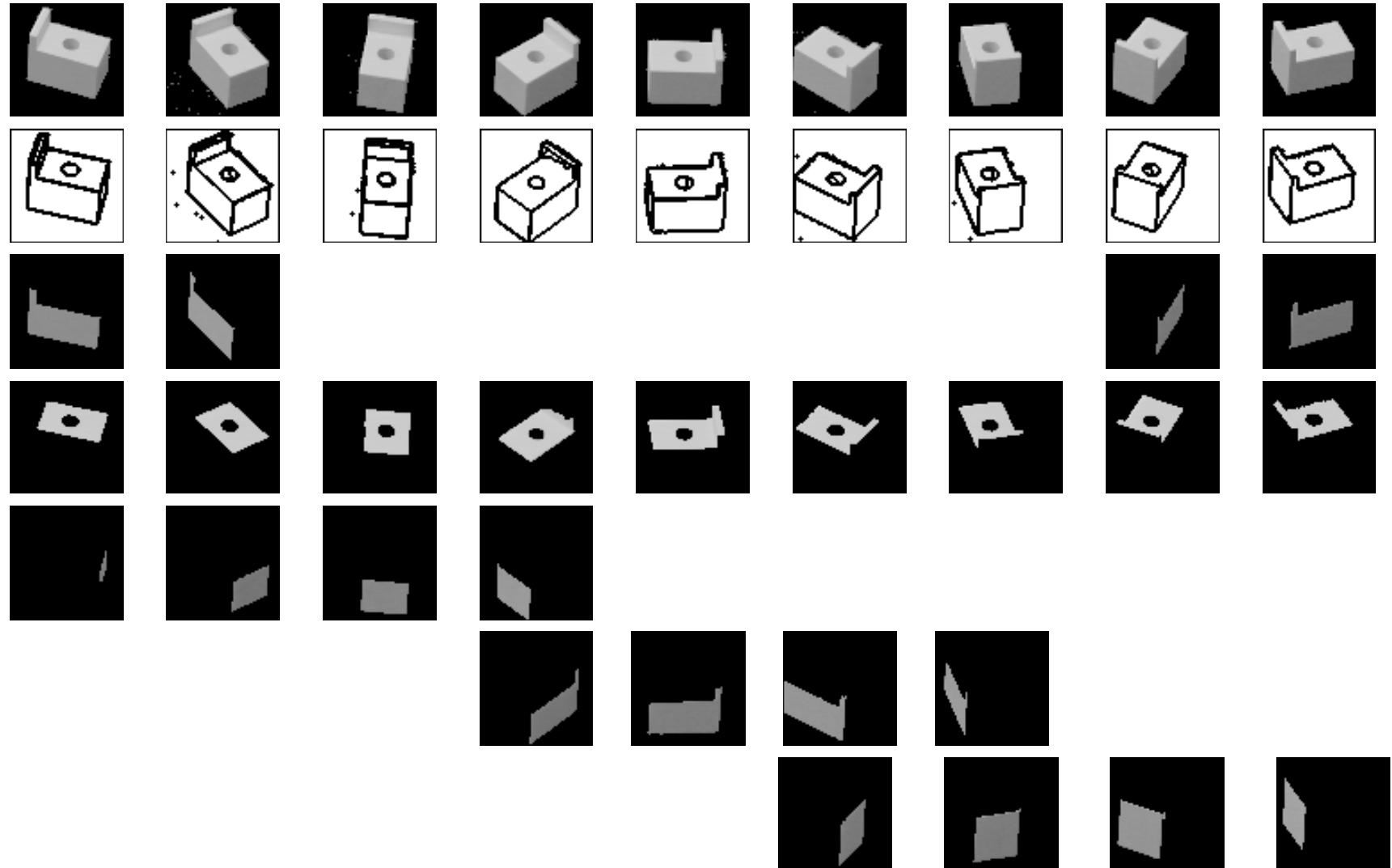
- Can model general-form objects
- Can deal with occlusion and clutter
- Can handle translation and scaling
- It is robust to segmentation problems

# Parts from Images

---

PARTS are image regions segmented using some segmentation algorithm.

# Appearance-Based Parts



# Appearance-Based Relations

---

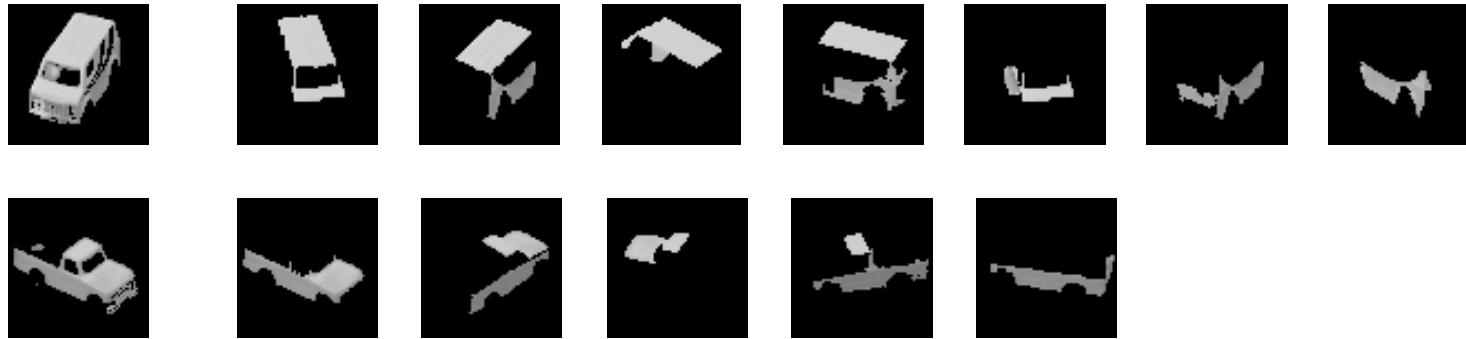
Several objects can have similar parts

Use spatial relations to discriminate

Represent relations using PCA

# Examples of ABRs

---



# Recognition in cluttered scenes

---

Segment the given image

PARTS:

- Project regions into the ABP eigenspace

- Make hypotheses if  $d < T_1$

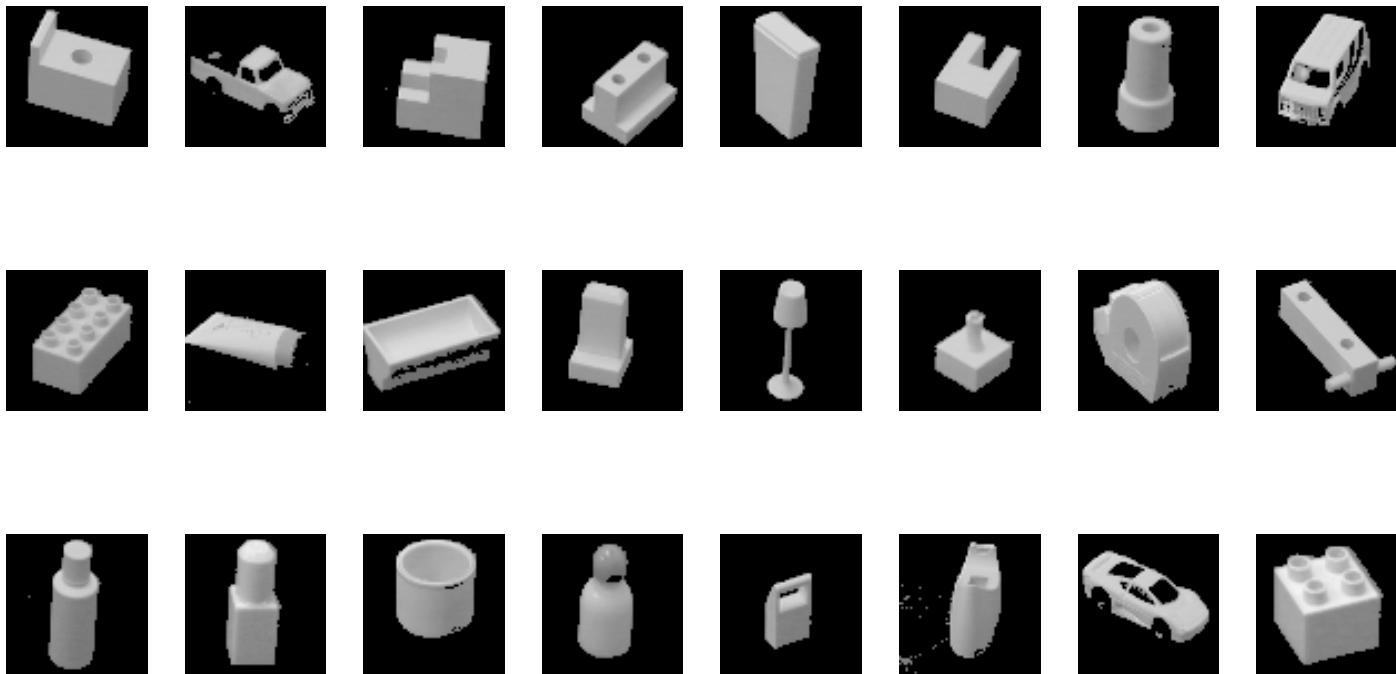
RELATIONS:

- Project adjacent regions with  $T_1 < d < T_2$  into the ABR eigenspace

- Make hypotheses if  $d' < T_3$

# Object Database

---

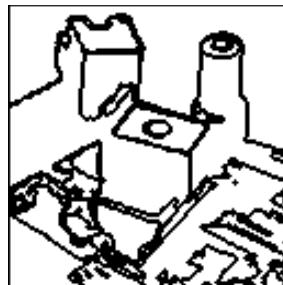
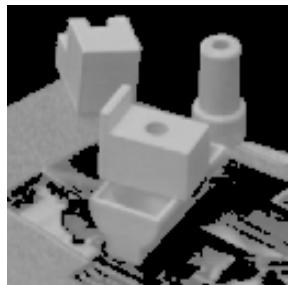


24 objects; 110 ABPs (69 shape groups/3 levels); 130 ABRs

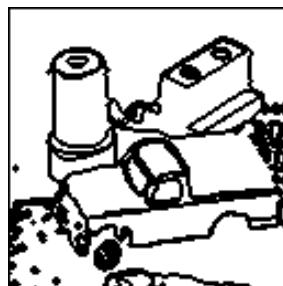
# Some Examples



<b>Identity</b>	<b>Score</b>	<b>Pose</b>
Sport Car	0.61	302.77
Truck	0.19	4.29
Ambulance	0.10	11.52



<b>Identity</b>	<b>Score</b>	<b>Pose</b>
Ccube	1.36	41.12
Stamp	0.89	0.64
Holecube	0.89	14.47
Sink	0.68	152.94



<b>Identity</b>	<b>Score</b>	<b>Pose</b>
Stamp	1.91	93.66
Twohole	0.96	181.63
Truck	0.4	93.67